

A gentle introduction to maximum entropy, log-linear, exponential, logistic, harmonic, Boltzmann, Markov Random Fields, Conditional Random Fields, etc., models

Mark Johnson

Department of Computing

March, 2013 (updated August 2015)

# A gentle introduction to maximum entropy, log-linear, exponential, logistic, harmonic, Boltzmann, Markov Random Field, etc., models

- How can we possibly cover so many kinds of models in a single talk?
- *Because they are all basically the same*
- If an idea is really (really!) good, you can justify it in many different ways!

# Outline

Introducing exponential models

Features in exponential models

Learning exponential models

Regularisation

Conditional models

Stochastic gradient descent and error-driven learning

Avoiding the partition function and its derivatives

Weakly labelled training data

Summary

# Why probabilistic models?

- Problem setup: given a set  $\mathcal{X}$  of possible items
  - ▶ e.g.,  $\mathcal{X}$  is the set of all possible English words (sequences of characters)
  - ▶ e.g.,  $\mathcal{X}$  is the set of all possible sentences (sequences of English words)
  - ▶ e.g.,  $\mathcal{X}$  is the set of all possible images ( $256 \times 256$  pixel arrays)
- Our goal is to learn a probability distribution  $\mathbf{P}(X)$  over  $\mathcal{X}$ 
  - ▶  $\mathbf{P}(X)$  identifies which items  $x \in \mathcal{X}$  are more likely and which ones are less likely
  - ▶ e.g., if  $\mathcal{X}$  is the set of possible English sentences,  $\mathbf{P}(X)$  is called a *language model*
  - ▶ language models are very useful in machine translation and speech recognition because they identify plausible sentences (e.g., “recognise speech” vs. “wreck a nice beach”)
- In this talk, we are interested in learning  $\mathbf{P}(X)$  from data  $D = (x_1, \dots, x_n)$ , which is sampled from the (unknown)  $\mathbf{P}(X)$

# Motivating exponential models

- **Goal:** define a probability distribution  $\mathbf{P}(X)$  over the  $x \in \mathcal{X}$
- Idea: describe  $x$  in terms of *weighted features*
- Let  $\mathcal{S}(x) \subseteq \mathcal{S}$  be the *set of  $x$ 's features*
- Let  $v_s$  be the *weight of feature  $s \in \mathcal{S}$* 
  - ▶ if  $v_s > 1$  then  $s$  makes  $x$  more probable
  - ▶ if  $v_s < 1$  then  $s$  makes  $x$  less probable
- If  $\mathcal{S}(x) = \{s_1, \dots, s_n\}$ , then

$$\begin{aligned}\mathbf{P}(X=x) &\propto v_{s_1} \dots v_{s_n} \\ &= \prod_{s \in \mathcal{S}(x)} v_s\end{aligned}$$

- Generalises many well-known models (e.g., HMMs, PCFGs)
  - ▶ what are the features and the feature weights in an HMM or a PCFG?



In generative models defined as a product of conditional distributions as factors, factors cannot be greater than 1

# The partition function

- Probability distributions must sum to 1, i.e.,

$$\sum_{x \in \mathcal{X}} \mathbf{P}(X=x) = 1$$

- But in general

$$\sum_{x \in \mathcal{X}} \left( \prod_{s \in \mathcal{S}(x)} v_s \right) \neq 1$$

⇒ *Normalise* the weighted feature products

$$Z = \sum_{x \in \mathcal{X}} \left( \prod_{s \in \mathcal{S}(x)} v_s \right)$$

$Z$  is called the *partition function*

- Then define:

$$\mathbf{P}(X=x) = \frac{1}{Z} \prod_{s \in \mathcal{S}(x)} v_s$$

Q: Why is  $Z$  called a *partition function*? What is it a function of?

# Feature functions

- Functions are often notationally easier to deal with than sets
- For each feature  $s \in \mathcal{S}$  define a *feature function*  $f_s : \mathcal{X} \mapsto 2$

$$f_s(x) = 1 \text{ if } s \in \mathcal{S}(x), \text{ and } 0 \text{ otherwise}$$

- Then we can rewrite

$$\begin{aligned} \mathbf{P}(X=x) &= \frac{1}{Z} \prod_{s \in \mathcal{S}(x)} v_s \\ &= \frac{1}{Z} \prod_{s \in \mathcal{S}} v_s^{f_s(x)} \end{aligned}$$

- Now we can have *real-valued feature functions*
- From here on assume we have a *vector of  $m$  feature functions*

$$\begin{aligned} \mathbf{f} &= (f_1, \dots, f_m), \text{ and} \\ \mathbf{f}(x) &= (f_1(x), \dots, f_m(x)) \end{aligned}$$

## Exponential form

- The feature weights  $v_j$  must be non-negative because probabilities are non-negative
- An easy way to ensure that feature weights are positive is to work in log space. Let  $w_j = \log(v_j)$  or equivalently  $v_j = \exp(w_j)$ .
  - ▶ If  $w_j > 0$  then having feature  $j$  makes  $x$  more probable
  - ▶ If  $w_j < 0$  then having feature  $j$  makes  $x$  less probable

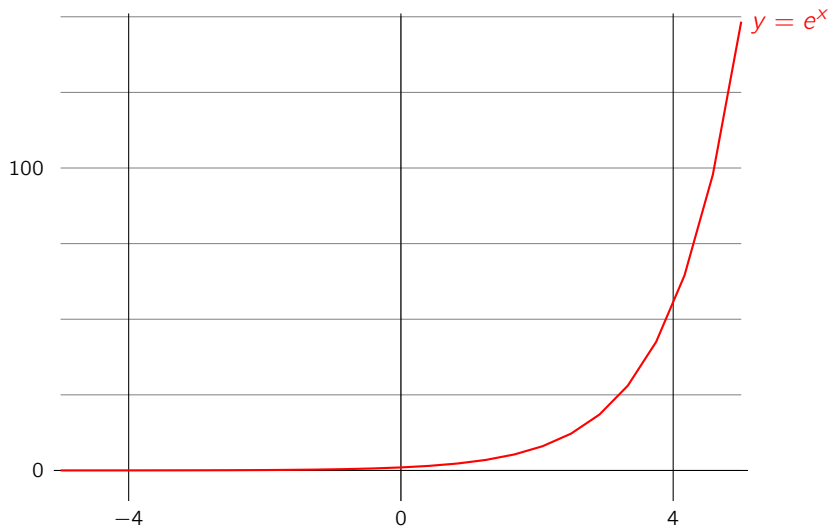
$$\begin{aligned}P(X=x) &= \frac{1}{Z} \prod_{j=1}^m v_j^{f_j(x)} \\ &= \frac{1}{Z} \exp\left(\sum_{j=1}^m w_j f_j(x)\right) \\ &= \frac{1}{Z} \exp(\mathbf{w} \cdot \mathbf{f}(x))\end{aligned}$$

where:

$$\begin{aligned}\mathbf{w} &= (w_1, \dots, w_m) \\ \mathbf{f}(x) &= (f_1(x), \dots, f_m(x)) \\ Z &= \sum_{x' \in \mathcal{X}} \exp(\mathbf{w} \cdot \mathbf{f}(x'))\end{aligned}$$



# The exponential function



# Outline

Introducing exponential models

Features in exponential models

Learning exponential models

Regularisation

Conditional models

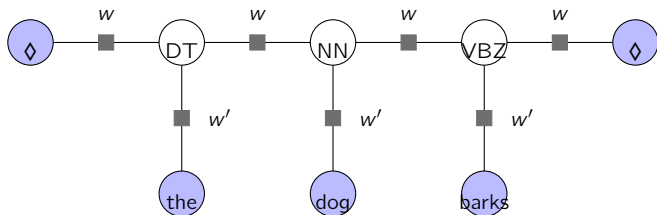
Stochastic gradient descent and error-driven learning

Avoiding the partition function and its derivatives

Weakly labelled training data

Summary

## Features in Random Fields



$$P(x) = \frac{1}{Z} w_{\diamond,DT} w'_{DT,the} w_{DT,NN} w'_{NN,dog} w_{NN,VBZ} w'_{VBZ,barks} w_{VBZ,\diamond}$$

- If  $\mathcal{V}$  is the set of words and  $\mathcal{Y}$  is the set of labels, there is a feature for each combination in  $\mathcal{Y} \times \mathcal{V}$  and for each combination in  $\mathcal{Y} \times \mathcal{Y}$ .
- If  $n_{y,y'}$  is the number of times label  $y$  precedes label  $y'$  in  $x$ , and  $m_{y,v}$  is the number of times label  $y$  appears with word  $v$ , then:

$$P(x) = \frac{1}{Z} \left( \prod_{y,y' \in \mathcal{Y} \times \mathcal{Y}} w_{y,y'}^{n_{y,y'}} \right) \left( \prod_{y,v \in \mathcal{Y} \times \mathcal{V}} w'_{y,v}^{m_{y,v}} \right)$$

# PCFGs and HMMs as exponential models

- Models like PCFGs and HMMs define the probability of a structure (e.g., a parse tree) *as a product of the probabilities of its components*
  - ▶ In a PCFG, each rule  $A \rightarrow \beta$  has a probability  $p_{A \rightarrow \beta}$
  - ▶ The probability of a tree is *the product of the probabilities of the rules used in its derivation*

$$P(x) = \prod_{A \rightarrow \beta \in \mathcal{R}} p_{A \rightarrow \beta}^{n_{A \rightarrow \beta}(x)}$$

where  $n_{A \rightarrow \beta}(x)$  is the number of times rule  $A \rightarrow \beta$  is used in derivation of tree  $x$

⇒ A PCFG can be expressed as an exponential model where:

- ▶ define a 1-to-1 mapping from PCFG rules to features (i.e., number the rules)
- ▶ define the feature functions:  $f_{A \rightarrow \beta}(x) = n_{A \rightarrow \beta}(x)$ , and
- ▶ set the feature values:  $v_{A \rightarrow \beta} = p_{A \rightarrow \beta}$

$$P(x) = \prod_{A \rightarrow \beta \in \mathcal{R}} v_{A \rightarrow \beta}^{f_{A \rightarrow \beta}(x)}$$

⇒ A PCFG (and an HMM) is an exponential model where  $Z = 1$

# Categorical features

- Suppose  $(g_1, \dots, g_m)$  are *categorical features*, where  $g_k$  ranges over  $\mathcal{G}_k$ 
  - ▶ E.g., if  $\mathcal{X}$  is a set of words, then  $\text{suffix}(x)$  might be the last letter of  $x$
- “One-hot” encoding of categorical features:
  - ▶ Define a binary feature  $f_{g_k=c}$  for each combination of a categorical feature  $g_k, k = 1, \dots, m$  and a possible value  $c \in \mathcal{G}_k$

$$f_{g_k=c}(x) = 1 \quad \text{if } g_k(x) = c, \text{ and } 0 \text{ otherwise}$$

- ⇒ Number of binary features grows extremely rapidly
- ▶ reranking parser has about 40 categorical features, but around 2 million binary features
  - But you *only need to instantiate feature-value pairs observed in training data*
    - ▶ learning procedures in general set  $w_{g=c} = 0$  if feature-value pair  $g(x) = c$  is not present in training data

## Feature redundancy in binary models

- Consider a situation where there are 2 outcomes:  $\mathcal{X} = \{a, b\}$

$$P(X=x) = \frac{1}{Z} \exp(\mathbf{w} \cdot \mathbf{f}(x)), \text{ where:}$$

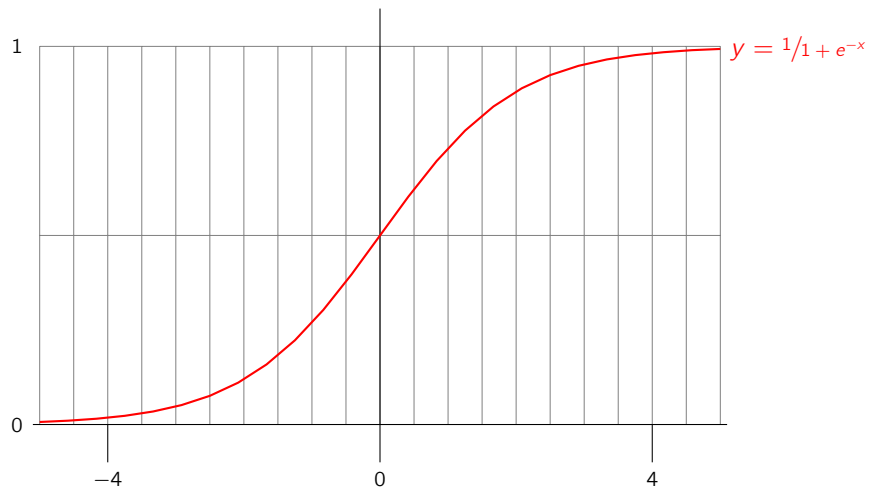
$$Z = \sum_{x' \in \mathcal{X}} \exp(\mathbf{w} \cdot \mathbf{f}(x')) = \exp(\mathbf{w} \cdot \mathbf{f}(a)) + \exp(\mathbf{w} \cdot \mathbf{f}(b)), \text{ so:}$$

$$\begin{aligned} P(X=a) &= \frac{\exp(\mathbf{w} \cdot \mathbf{f}(a))}{\exp(\mathbf{w} \cdot \mathbf{f}(a)) + \exp(\mathbf{w} \cdot \mathbf{f}(b))} \\ &= \frac{1}{1 + \exp(\mathbf{w} \cdot (\mathbf{f}(b) - \mathbf{f}(a)))} \\ &= s(\mathbf{w} \cdot (\mathbf{f}(a) - \mathbf{f}(b))), \text{ where:} \end{aligned}$$

$$s(z) = \frac{1}{1 + \exp(-z)} \text{ is the } \textit{logistic sigmoid function}$$

⇒ In binary models *only the difference between feature values matters*

# The logistic sigmoid function



## Feature redundancy in exponential models

- This result generalises to all exponential models
- Let  $\mathbf{u} = (u_1, \dots, u_m)$  be *any vector* of same dimensionality as the features
- Define an exponential model using *new feature functions*  $\mathbf{f}'(x) = \mathbf{f}(x) + \mathbf{u}$ .  
Then:

$$\begin{aligned} P(X=x) &= \frac{\exp(\mathbf{w} \cdot \mathbf{f}'(x))}{\sum_{x' \in \mathcal{X}} \exp(\mathbf{w} \cdot \mathbf{f}'(x'))} \\ &= \frac{\exp(\mathbf{w} \cdot \mathbf{f}(x)) \exp(\mathbf{w} \cdot \mathbf{u})}{\sum_{x' \in \mathcal{X}} \exp(\mathbf{w} \cdot \mathbf{f}(x')) \exp(\mathbf{w} \cdot \mathbf{u})} \\ &= \frac{\exp(\mathbf{w} \cdot \mathbf{f}(x))}{\sum_{x' \in \mathcal{X}} \exp(\mathbf{w} \cdot \mathbf{f}(x'))} \end{aligned}$$

$\Rightarrow$  *Adding or subtracting a constant vector to feature values does not change the distribution defined by an exponential model*

- The feature extractor for the reranking parser subtracts the vector  $\mathbf{u}$  that *makes the feature vectors as sparse as possible*



# Outline

Introducing exponential models

Features in exponential models

**Learning exponential models**

Regularisation

Conditional models

Stochastic gradient descent and error-driven learning

Avoiding the partition function and its derivatives

Weakly labelled training data

Summary

# Methods for learning from data

- Learning or estimating feature weights  $\mathbf{w}$  from *training data*  $D = (x_1, \dots, x_n)$ , where each  $x_j \in \mathcal{X}$
- *Maximum likelihood*: choose  $\mathbf{w}$  to make  $D$  as likely as possible

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmax}} L_D(\mathbf{w}), \text{ where:}$$

$$L_D(\mathbf{w}) = \prod_{i=1}^n P_{\mathbf{w}}(x_i)$$

- *Minimising negative log likelihood* is mathematically equivalent, and has mathematical and computational advantages
  - ▶ negative log likelihood is *convex* (with fully visible training data)
  - ⇒ single optimum that can be found by “following gradient downhill”
  - ▶ avoids floating point underflow
- But other learning methods may have advantages
  - ▶ with a large number of features, a *regularisation penalty term* (e.g., L1 and/or L2 prior) helps to avoid overfitting
  - ▶ *optimising a specialised loss function* (e.g., expected f-score) can improve performance on a specific task

# Learning as minimising a loss function

- Goal: find the feature weights  $\hat{\mathbf{w}}$  that *minimise the negative log likelihood*  $\ell_D$  of feature weights  $\mathbf{w}$  given data  $D = (x_1, \dots, x_n)$ :

$$\begin{aligned}\hat{\mathbf{w}} &= \underset{\mathbf{w}}{\operatorname{argmin}} \ell_D(\mathbf{w}) \\ \ell_D(\mathbf{w}) &= -\log L_D(\mathbf{w}) = -\log \prod_{i=1}^n \mathbf{P}_{\mathbf{w}}(x_i) \\ &= \sum_{i=1}^n -\log \mathbf{P}_{\mathbf{w}}(x_i)\end{aligned}$$

- The negative log likelihood  $\ell_D$  is a sum of the *losses*  $-\log \mathbf{P}_{\mathbf{w}}(x_i)$  the model  $\mathbf{w}$  incurs on each data item  $x_i$
- The maximum likelihood estimator selects the model  $\hat{\mathbf{w}}$  that minimises the loss  $\ell_D$  on data set  $D$
- Many other machine learning algorithms for estimating  $\mathbf{w}$  from  $D$  can be understood as minimising some loss function

# Why is learning exponential models hard?

- Exponential models are so flexible because the features can have arbitrary weights
- ⇒ The partition function  $Z$  is required to ensure the distribution  $\mathbf{P}(x)$  is normalised
- The partition function  $Z$  *varies as a function of  $\mathbf{w}$*

$$\mathbf{P}(X=x) = \frac{1}{Z} \exp(\mathbf{w} \cdot \mathbf{f}(x)), \text{ where:}$$
$$Z = \sum_{x' \in \mathcal{X}} \exp(\mathbf{w} \cdot \mathbf{f}(x'))$$

- ⇒ So we can't ignore  $Z$ , which makes it hard to optimise the likelihood!
  - ▶ no closed-form solution for the feature weights  $w_j$
  - ▶ learning usually involves *numerically optimising* the likelihood function or some other loss function
  - ▶ calculating  $Z$  requires *summing over entire space  $\mathcal{X}$*
  - ▶ many methods for approximating  $Z$  and/or its derivatives; typically unclear how the approximations affect the estimates of  $\mathbf{w}$

# The derivative of the negative log likelihood

- Efficient numerical optimisation routines require evaluation of the function to be minimised (negative log likelihood  $\ell_D$ ) *and its derivatives*
  - ▶ use a standard package; L-BFGS (LMVM), conjugate gradient
- We'll optimise  $1/n$  times the negative log likelihood of  $\mathbf{w}$  given data  $D = (x_1, \dots, x_n)$ :

$$\ell_D(\mathbf{w}) = -\frac{1}{n} \sum_{i=1}^n \log \mathbf{P}_{\mathbf{w}}(x_i) = \log Z - \frac{1}{n} \sum_{i=1}^n \mathbf{w} \cdot \mathbf{f}(x_i)$$

- The derivative of  $\ell$  is:

$$\frac{\partial \ell_D}{\partial w_j} = \mathbf{E}_{\mathbf{w}}[f_j] - \mathbf{E}_D[f_j], \text{ where:}$$

$$\mathbf{E}_{\mathbf{w}}[f_j] = \sum_{x' \in \mathcal{X}} f_j(x') \mathbf{P}_{\mathbf{w}}(x') \quad (\text{expected value of } f_j \text{ wrt } \mathbf{P}_{\mathbf{w}})$$

$$\mathbf{E}_D[f_j] = \frac{1}{n} \sum_{i=1}^n f_j(x_i) \quad (\text{expected value of } f_j \text{ wrt } D)$$

- At optimum  $\partial \ell_D / \partial \mathbf{w} = \mathbf{0}$   
 $\Rightarrow$  *model's expected feature values equals data's feature values*

## Exercise: derive the formulae on the previous slide!

- This is a basic result for exponential models that is the basis of many other results
- If you want to generalise exponential models, you'll need to derive similar formulae
- You'll need to know:
  - ▶ that derivatives distribute over sums
  - ▶ that  $\partial \log(x) / \partial x = 1/x$
  - ▶ the chain rule, i.e., that  $\partial y / \partial x = \partial y / \partial u \partial u / \partial x$

# Maximum entropy models

- Idea: given training data  $D$  and feature functions  $\mathbf{f}$ , find the distribution  $\mathbf{P}'(X)$  that:
  1.  $\mathbf{E}_{\mathbf{P}'}[f_j] = \mathbf{E}_D[f_j]$  for all features  $f_j$ ,  
i.e.,  $\mathbf{P}'$  agrees with  $D$  on the features
  2. of all distributions satisfying (1),  $\mathbf{P}'$  has *maximum entropy*  
i.e.,  $\mathbf{P}'$  has the least possible additional information
- Because  $\widehat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} \ell_D(\mathbf{w})$  then

$$\frac{\partial \ell_D}{\partial \mathbf{w}}(\widehat{\mathbf{w}}) = \mathbf{0}$$

$\Rightarrow \mathbf{E}_{\mathbf{w}}[f_j] = \mathbf{E}_D[f_j]$  for all features  $f_j$

- Theorem:  $\mathbf{P}_{\mathbf{w}} = \mathbf{P}'$ , i.e., for any data  $D$  and feature functions  $\mathbf{f}$  *the maximum likelihood distribution and the maximum entropy distribution are the same distribution*

# Outline

Introducing exponential models

Features in exponential models

Learning exponential models

**Regularisation**

Conditional models

Stochastic gradient descent and error-driven learning

Avoiding the partition function and its derivatives

Weakly labelled training data

Summary



# Why regularise?

- If every  $x \in D$  has feature  $f_j$  and some  $x \in \mathcal{X}$  does not, then  $\widehat{w}_j = \infty$
- If no  $x \in D$  has feature  $f_j$  and some  $x \in \mathcal{X}$  does, then  $\widehat{w}_j = -\infty$
- Infinities cause problems for numerical routines
- Just because a feature always occurs/doesn't occur in training data doesn't mean this will also occur in test data ("accidental zeros")
- These are extreme examples of *overlearning*
  - ▶ overlearning often occurs when the size of the data  $D$  is not much greater than the number of features  $m$
- Idea: add a *regulariser* (also called a penalty term or prior) to the negative log likelihood that *penalises large feature weights*
  - ▶ Recall that  $w_j = 0$  means that feature  $f_j$  is ignored

## $L_2$ regularisation

- Instead of minimising the negative log likelihood  $\ell_D(\mathbf{w})$ , we optimise

$$\begin{aligned}\widehat{\mathbf{w}} &= \underset{\mathbf{w}}{\operatorname{argmin}} \ell_D(\mathbf{w}) + c R(\mathbf{w}), \text{ where:} \\ R(\mathbf{w}) &= \|\mathbf{w}\|_2^2 \\ &= \mathbf{w} \cdot \mathbf{w} \\ &= \sum_{j=1}^m w_j^2\end{aligned}$$

- $R$  is a *penalty term* that varies with each feature weight  $w_j$  such that:
  - ▶ the penalty is zero when  $w_j = 0$ ,
  - ▶ the penalty is *greater than zero whenever  $w_j \neq 0$* , and
  - ▶ the penalty *grows as  $w_j$  moves further away from 0*
- The *regulariser constant*  $c$  is usually set to *optimise performance on held-out data*

# Bayesian MAP estimation

- Recall Bayesian belief updating:

$$\underbrace{P(\text{Hypothesis} \mid \text{Data})}_{\text{Posterior}} \propto \underbrace{P(\text{Data} \mid \text{Hypothesis})}_{\text{Likelihood}} \underbrace{P(\text{Hypothesis})}_{\text{Prior}}$$

- In our setting:
  - ▶ **Data** =  $D = (x_1, \dots, x_n)$
  - ▶ **Hypothesis** =  $\mathbf{w} = (w_1, \dots, w_m)$
- If we want the MAP (Maximum Aposteriori) estimate for  $\mathbf{w}$ :

$$\begin{aligned}\hat{\mathbf{w}} &= \operatorname{argmax}_{\mathbf{w}} \underbrace{P(\mathbf{w} \mid D)}_{\text{Posterior}} \\ &= \operatorname{argmax}_{\mathbf{w}} \underbrace{P(D \mid \mathbf{w})}_{\text{Likelihood}} \underbrace{P(\mathbf{w})}_{\text{Prior}} \\ &= \operatorname{argmax}_{\mathbf{w}} \left( \prod_{i=1}^n P(x_i \mid \mathbf{w}) \right) P(\mathbf{w})\end{aligned}$$

## Regularisation as Bayesian MAP estimation

- Restate the MAP estimate in terms of negative log likelihood  $\ell_D$ :

$$\begin{aligned}\hat{\mathbf{w}} &= \underset{\mathbf{w}}{\operatorname{argmax}} \underbrace{\left( \prod_{i=1}^n \mathbf{P}(x_i | \mathbf{w}) \right)}_{\text{Likelihood}} \underbrace{\mathbf{P}(\mathbf{w})}_{\text{Prior}} \\ &= \underset{\mathbf{w}}{\operatorname{argmin}} \left( - \sum_{i=1}^n \log \mathbf{P}(x_i | \mathbf{w}) \right) - \log \mathbf{P}(\mathbf{w}) \\ &= \underset{\mathbf{w}}{\operatorname{argmin}} \ell_D(\mathbf{w}) - \log \mathbf{P}(\mathbf{w}), \text{ where:}\end{aligned}$$

$$\ell_D(\mathbf{w}) = - \sum_{i=1}^n \log \mathbf{P}(x_i | \mathbf{w})$$

$\Rightarrow$  MAP estimate  $\hat{\mathbf{w}}$  equals regularised MLE  $\hat{\mathbf{w}}$

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \ell_D(\mathbf{w}) + c R(\mathbf{w})$$

if  $cR(\mathbf{w}) = -\log \mathbf{P}(\mathbf{w})$ , i.e., *if the regulariser is the negative log prior*

## $L_2$ regularisation as a Gaussian prior

- What kind of prior is an  $L_2$  regulariser?
- If  $cR(\mathbf{w}) = -\log \mathbf{P}(\mathbf{w})$  then

$$\mathbf{P}(\mathbf{w}) = \exp(-cR(\mathbf{w}))$$

- If  $R(\mathbf{w}) = \|\mathbf{w}\|_2^2 = \sum_{j=1}^m w_j^2$ , then *the prior is a zero-mean Gaussian*

$$\mathbf{P}(\mathbf{w}) \propto \exp\left(-c \sum_{j=1}^m w_j^2\right)$$

The additional factors in the Gaussian become constants in log probability space, and therefore can be ignored when finding  $\hat{\mathbf{w}}$

- $L_2$  regularisation is also known as *ridge regularisation*

# $L_1$ regularisation or Lasso regularisation

- The  $L_1$  norm is the sum of the absolute values

$$\begin{aligned}R(\mathbf{w}) &= \|\mathbf{w}\|_1 \\ &= \sum_{j=1}^m |w_j|\end{aligned}$$

- $L_1$  regularisation is popular because it produces *sparse feature weights*
  - ▶ a feature weight vector  $\mathbf{w}$  is *sparse* iff most of its values are zero
- But it's difficult to optimise  $L_1$ -regularised log-likelihood because *its derivative is discontinuous at the orthant boundaries*

$$\frac{\partial R}{\partial w_j} = \begin{cases} +1 & \text{if } w_j > 0 \\ -1 & \text{if } w_j < 0 \end{cases}$$

- Specialised versions of standard numerical optimisers have been developed to optimise  $L_1$ -regularised log-likelihood

# What does regularisation do?

- Regularised negative log likelihood

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \ell_D(\mathbf{w}) + cR(\mathbf{w})$$

- At the optimum weights  $\hat{\mathbf{w}}$ , for each  $j$ :

$$\frac{\partial \ell_D}{\partial w_j} + c \frac{\partial R}{\partial w_j} = 0, \text{ or equivalently}$$

$$\mathbf{E}_D[f_j] - \mathbf{E}_{\mathbf{w}}[f_j] = c \frac{\partial R}{\partial w_j}$$

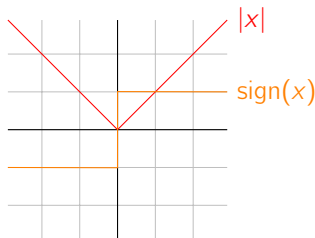
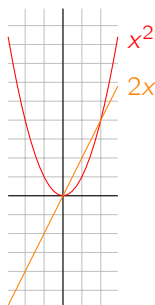
I.e., the regulariser gives the model some “slack” in requiring the empirical expected feature values equal the model’s predicted expected feature values.

# Why does $L_1$ regularisation produce sparse weights?

- Regulariser's derivative specifies gap between empirical and model feature expectation

$$\mathbf{E}_D[f_j] - \mathbf{E}_w[f_j] = c \frac{\partial R}{\partial w_j}$$

- For  $L_2$  regularisation,  
 $\frac{\partial R}{\partial w_j} \rightarrow 0$  as  $w_j \rightarrow 0$ 
  - ▶ little effect on small  $w$ $\Rightarrow$  no reason for feature weights to be zero
- For  $L_1$  regularisation,  
 $\frac{\partial R}{\partial w_j} \rightarrow \text{sign}(w_j)$  as  $w_j \rightarrow 0$ 
  - ▶ regulariser has effect whenever  $w \neq 0$
  - ▶ regulariser drives feature weights to 0 whenever “expectation gap”  $< c$





## Group sparsity via the Group Lasso

- Sometimes features come in natural groups; e.g.,  $\mathbf{F} = (\mathbf{f}_1, \dots, \mathbf{f}_m)$ , where each  $\mathbf{f}_j = (f_{j,1}, \dots, f_{j,v_j}), j = 1, \dots, m$
- Corresponding weights  $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_m)$ , where each  $\mathbf{w}_j = (w_{j,1}, \dots, w_{j,v_j})$

$$\mathbf{P}(X=x) = \frac{1}{Z} \exp \left( \sum_{j=1}^m \sum_{k=1}^{v_j} w_{j,k} f_{j,k}(x) \right)$$

- We'd like *group sparsity*, i.e., for “most”  $j \in 1, \dots, m$ ,  $\mathbf{w}_j = \mathbf{0}$
- The *group Lasso* regulariser achieves this:

$$\begin{aligned} R(\mathbf{W}) &= \sum_{j=1}^m c_j \|\mathbf{w}_j\|_2 \\ &= \sum_{j=1}^m c_j \left( \sum_{k=1}^{v_j} w_{j,k}^2 \right)^{1/2} \end{aligned}$$

# Optimising the regularised log likelihood

- Learning feature weights involves optimising regularised likelihood

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \ell_D(\mathbf{w}) + cR(\mathbf{w})$$

$$\ell_D(\mathbf{w}) = -\frac{1}{n} \sum_{i=1}^n \log \mathbf{P}(x_i) = \log Z - \frac{1}{n} \sum_{i=1}^n \mathbf{w} \cdot \mathbf{f}(x_i)$$

$$Z = \sum_{x' \in \mathcal{X}} \exp(\mathbf{w} \cdot \mathbf{f}(x'))$$

- Challenges in optimisation:
    - ▶ If regulariser  $R$  is not differentiable (e.g.,  $R = L_1$ ), then you need a specialised optimisation algorithm to handle discontinuous derivatives
    - ▶ if  $\mathcal{X}$  is large (infinite), calculating  $Z$  may be difficult because it involves summing over  $\mathcal{X}$
- ⇒ just evaluate on the subset  $\mathcal{X}' \subset \mathcal{X}$  where  $\mathbf{w} \cdot \mathbf{f}$  is largest (assuming you can find it)

# Outline

Introducing exponential models

Features in exponential models

Learning exponential models

Regularisation

**Conditional models**

Stochastic gradient descent and error-driven learning

Avoiding the partition function and its derivatives

Weakly labelled training data

Summary

# Why conditional models?

- In a conditional model, each datum is a pair  $(x, y)$ , where  $x \in \mathcal{X}$  and  $y \in \mathcal{Y}$
- The goal of a conditional model is to *predict  $y$  given  $x$*
- Usually  $x$  is an item or an observation and  $y$  is a *label* for  $x$ 
  - ▶ e.g.,  $\mathcal{X}$  is the set of all possible news articles, and  $\mathcal{Y}$  is a set of topics, e.g.  $\mathcal{Y} = \{\text{finance, sports, politics, } \dots\}$
  - ▶ e.g.,  $\mathcal{X}$  is the set of all possible  $256 \times 256$  images, and  $\mathcal{Y}$  is a set of labels, e.g.,  $\mathcal{Y} = \{\text{cat, dog, person, } \dots\}$
  - ▶ e.g.,  $\mathcal{X}$  is the set of all possible Tweets, and  $\mathcal{Y}$  is a Boolean value indicating whether  $x \in \mathcal{X}$  expresses a sentiment
  - ▶ e.g.,  $\mathcal{X}$  is the set of all possible sentiment-expressing Tweets, and  $\mathcal{Y}$  is a Boolean value indicating whether  $x \in \mathcal{X}$  has positive or negative sentiment
- We will do this by learning a *conditional probability distribution*  $\mathbf{P}(Y | X)$ , which is the probability of  $Y$  *given*  $X$
- We estimate  $\mathbf{P}(Y | X)$  from data  $D = ((x_1, y_1), \dots, (x_n, y_n))$ , that consists of pairs of items  $x_i$  and their labels  $y_i$  (*supervised learning*)
  - ▶ in *unsupervised learning* we are only given the data items  $x_i$ , but not their labels  $y_i$  (clustering)
  - ▶ in *semi-supervised learning* we are not given the labels  $y_i$  for all data items  $x_i$  (we might be given only some labels, or the labels might only be partially identified)

# Conditional exponential models

- Data  $D = ((x_1, y_1), \dots, (x_n, y_n))$  consists of  $(x, y)$  *pairs*, where  $x \in \mathcal{X}$  and  $y \in \mathcal{Y}$
- Want to predict  $y$  from  $x$ , for which we only need *conditional distribution*  $\mathbf{P}(Y | X)$ , not the joint distribution  $\mathbf{P}(X, Y)$
- Features are now functions  $f(x, y)$  over  $(x, y)$  pairs
- Conditional exponential model:

$$\mathbf{P}(y | x) = \frac{1}{Z(x)} \exp(\mathbf{w} \cdot \mathbf{f}(x, y)), \text{ where:}$$
$$Z(x) = \sum_{y' \in \mathcal{Y}} \exp(\mathbf{w} \cdot \mathbf{f}(x, y'))$$

- Big advantage:  $Z(x)$  only requires a sum over  $\mathcal{Y}$ , while “joint” partition function  $Z$  requires a sum over all  $\mathcal{X} \times \mathcal{Y}$  pairs
  - ▶ in many applications label set  $\mathcal{Y}$  is small
  - ▶ *size of  $\mathcal{X}$  doesn't affect computational effort to compute  $Z(x)$*

# Features in conditional models

- In a conditional model, changing the feature function  $\mathbf{f}(x, y)$  to  $\mathbf{f}'(x, y) = \mathbf{f}(x, y) + \mathbf{u}(x)$  *does not change the distribution*  $\mathbf{P}(y | x)$ 
  - ⇒ *adding or subtracting a function that only depends on  $x$  does not affect a conditional model*
  - ⇒ *to be useful in a conditional model, a feature must be a non-constant function of  $y$*
- A feature  $f(x, y) = f(y)$  that only depends on  $y$  behaves like a *bias node* in a neural net
  - ▶ it's often a good idea to have a “one-hot” feature for each  $c \in \mathcal{Y}$ :

$$f_{y=c}(y) = 1 \text{ if } y = c, \text{ and } 0 \text{ otherwise}$$

- If  $\mathcal{X}$  is a set of discrete categories, it's often useful to have pairwise “one-hot” features for each  $c \in \mathcal{X}$  and  $c' \in \mathcal{Y}$

$$f_{x=c, y=c'}(x, y) = 1 \text{ if } x = c \text{ and } y = c', \text{ and } 0 \text{ otherwise}$$

# Using a conditional model to make predictions

- Labelling problem: we have feature weights  $\mathbf{w}$  and want to predict label  $y$  for some  $x$
- The most probable label  $\hat{y}(x)$  given  $x$  is:

$$\begin{aligned}\hat{y}(x) &= \operatorname{argmax}_{y' \in \mathcal{Y}} \mathbf{P}_{\mathbf{w}}(Y=y' \mid X=x) \\ &= \operatorname{argmax}_{y' \in \mathcal{Y}} \frac{1}{Z(x)} \exp(\mathbf{w} \cdot \mathbf{f}(x, y')) \\ &= \operatorname{argmax}_{y' \in \mathcal{Y}} \mathbf{w} \cdot \mathbf{f}(x, y')\end{aligned}$$

- Partition function  $Z(x)$  is a constant here, so drops out

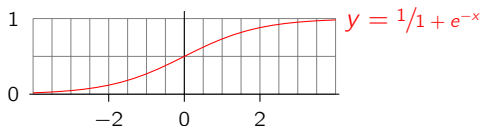
# Logistic regression

- Suppose  $\mathcal{Y} = \{0, 1\}$ , i.e., our labels are Boolean

$$\begin{aligned} P(Y=1 | X=x) &= \frac{\exp(\mathbf{w} \cdot \mathbf{f}(x, 1))}{\exp(\mathbf{w} \cdot \mathbf{f}(x, 0)) + \exp(\mathbf{w} \cdot \mathbf{f}(x, 1))} \\ &= \frac{1}{1 + \exp(\mathbf{w} \cdot (\mathbf{f}(x, 0) - \mathbf{f}(x, 1)))} \\ &= \frac{1}{1 + \exp(-\mathbf{w} \cdot \mathbf{g}(x))}, \text{ where:} \\ \mathbf{g}(x) &= \mathbf{f}_j(x, 1) - \mathbf{f}_j(x, 0), \text{ for all } j \in 1, \dots, m \end{aligned}$$

⇒ Only *relative feature differences matter* in a conditional model

- Logistic sigmoid function:





# Estimating conditional exponential models

- Compute *maximum conditional likelihood estimator* by minimizing negative log conditional likelihood

$$\begin{aligned}\hat{\mathbf{w}} &= \underset{\mathbf{w}}{\operatorname{argmin}} \ell_D(\mathbf{w}), \text{ where:} \\ \ell_D(\mathbf{w}) &= -\sum_{i=1}^n \log \mathbf{P}_{\mathbf{w}}(y_i | x_i) \\ &= \sum_{i=1}^n (\log Z(x_i) - \mathbf{w} \cdot \mathbf{f}(x_i, y_i))\end{aligned}$$

- Derivatives are differences of *conditional expectations* and *empirical feature values*

$$\begin{aligned}\frac{\partial \ell_D}{\partial w_j} &= \sum_{i=1}^n (\mathbf{E}_{\mathbf{w}}[f_j | x_i] - f_j(x_i, y_i)), \text{ where:} \\ \mathbf{E}_{\mathbf{w}}[f_j | x] &= \sum_{y' \in \mathcal{Y}} f_j(x, y') \mathbf{P}_{\mathbf{w}}(y' | x) \quad (\text{expected value of } f_j \text{ given } x)\end{aligned}$$

# Regularising conditional exponential models

- Calculating derivatives of conditional likelihood only requires summing over  $\mathcal{Y}$ , and not  $\mathcal{X}$ 
  - ▶ not too expensive if  $|\mathcal{Y}|$  is small
  - ▶ if  $\mathcal{Y}$  has a regular structure (e.g., a sequence), then there may be efficient algorithms for summing over  $\mathcal{Y}$
- Regularisation adds a *penalty term* to objective function we seek to optimise
- *Important to regularise* (unless number of features is small)
  - ▶  $L_1$  (Lasso) regularisation produces *sparse feature weights*
  - ▶  $L_2$  (ridge) regularisation produces dense feature weights
  - ▶ Group lasso regularisation produces group-level sparsity in feature weights

# Outline

Introducing exponential models

Features in exponential models

Learning exponential models

Regularisation

Conditional models

Stochastic gradient descent and error-driven learning

Avoiding the partition function and its derivatives

Weakly labelled training data

Summary

# Why stochastic gradient descent?

- For small/medium data sets, “batch” methods using standard numerical optimisation procedures (such as L-BFGS) can work very well
  - ▶ these directly *minimise the negative log likelihood*  $\ell_D$
  - ▶ to calculate the negative log likelihood and its derivatives requires a pass through the entire training data
- But for very large data sets (e.g., data sets that don't fit into memory), or with very large models (such as neural nets), these can be too slow
- *Stochastic gradient descent* calculates a noisy gradient from a small subset of the training data, so it can learn considerably faster
  - ▶ but the solution it finds is often less accurate

# Gradient descent and mini-batch algorithms

- Idea: to minimise  $\ell_D(\mathbf{w})$ , move in direction of negative gradient  $\partial\ell_D/\partial\mathbf{w}$
- If  $\widehat{\mathbf{w}}^{(t)}$  is current estimate of  $\mathbf{w}$ , update as follows:

$$\begin{aligned}\widehat{\mathbf{w}}^{(t+1)} &= \widehat{\mathbf{w}}^{(t)} - \varepsilon \frac{\partial\ell_D}{\partial\mathbf{w}}(\widehat{\mathbf{w}}^{(t)}) \\ &= \widehat{\mathbf{w}}^{(t)} - \varepsilon \sum_{i=1}^n (\mathbb{E}_{\widehat{\mathbf{w}}^{(t)}}[\mathbf{f} \mid x_i] - \mathbf{f}(x_i, y_i))\end{aligned}$$

- $\varepsilon$  is *step size*; can be difficult to find a good value for it!
- This is *not a good optimisation algorithm*, as it zig-zags across valleys
- Update is *difference between expected and empirical feature values*
  - ▶ Each update requires a full pass through  $D \Rightarrow$  relatively slow
- “Mini-batch algorithms”: calculate expectations on *a small sample of  $D$*  to determine weight updates

# Stochastic gradient descent as mini-batch of size 1

- Stochastic Gradient Descent (SGD) is the mini-batch algorithm with a mini-batch of size 1
- If  $\widehat{\mathbf{w}}^{(t)}$  is current estimate of  $\mathbf{w}$ , training data  $D = ((x_1, y_1), \dots, (x_n, y_n))$ , and  $r_t$  is a random number in  $1, \dots, n$  then:

$$\widehat{\mathbf{w}}^{(t+1)} = \widehat{\mathbf{w}}^{(t)} - \varepsilon (\mathbf{E}_{\widehat{\mathbf{w}}^{(t)}}[\mathbf{f} \mid x_{r_t}] - \mathbf{f}(x_{r_t}, y_{r_t})), \text{ where:}$$

$$\mathbf{E}_{\mathbf{w}}[\mathbf{f} \mid x] = \sum_{y' \in \mathcal{Y}} \mathbf{f}(x, y') \mathbf{P}_{\mathbf{w}}(y' \mid x)$$

$$\mathbf{P}_{\mathbf{w}}(y \mid x) = \frac{1}{Z(x)} \exp(\mathbf{w} \cdot \mathbf{f}(x, y))$$

$$Z(x) = \sum_{y' \in \mathcal{Y}} \exp(\mathbf{w} \cdot \mathbf{f}(x, y'))$$

- Stochastic gradient descent updates estimate of  $\mathbf{w}$  after seeing each training example
- ⇒ Learning can be very fast; might not even need a full pass over  $D$
- Perhaps the most widely used learning algorithm today

# The Perceptron algorithm as approximate SGD

- Idea: assume  $\mathbf{P}_{\mathbf{w}}(y | x)$  is peaked around  $\hat{y}_{\mathbf{w}}(x) = \operatorname{argmax}_{y' \in \mathcal{Y}} \mathbf{w} \cdot \mathbf{f}(x, y')$ .  
Then:

$$\begin{aligned} \mathbf{E}_{\mathbf{w}}[\mathbf{f} | x] &= \sum_{y' \in \mathcal{Y}} \mathbf{f}(x, y') \mathbf{P}_{\mathbf{w}}(y' | x) \\ &\approx \mathbf{f}(x, \hat{y}_{\mathbf{w}}(x)) \end{aligned}$$

- Plugging this into the SGD algorithm, we get:

$$\begin{aligned} \hat{\mathbf{w}}^{(t+1)} &= \hat{\mathbf{w}}^{(t)} - \varepsilon \left( \mathbf{E}_{\hat{\mathbf{w}}^{(t)}}[\mathbf{f} | x_{r_t}] - \mathbf{f}(x_{r_t}, y_{r_t}) \right) \\ &\approx \hat{\mathbf{w}}^{(t)} - \varepsilon \left( \mathbf{f}(x_{r_t}, \hat{y}_{\hat{\mathbf{w}}^{(t)}}(x)) - \mathbf{f}(x_{r_t}, y_{r_t}) \right) \end{aligned}$$

- This is an *error-driven learning rule*, since *no update is made on iteration  $t$  if  $\hat{y}(x_{r_t}) = y_{r_t}$*

# Regularisation as weight decay in SGD and Perceptron

- Regularisation: minimise a *penalised negative log likelihood*

$$\widehat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \ell_D(\mathbf{w}) + cR(\mathbf{w}), \text{ where:}$$
$$R(\mathbf{w}) = \begin{cases} \sum_{j=1}^m w_j^2 & \text{with an } L_2 \text{ regulariser} \\ \sum_{j=1}^m |w_j| & \text{with an } L_1 \text{ regulariser} \end{cases}$$

- Adding  $L_2$  regularisation in SGD and Perceptron introduces *multiplicative weight decay*:

$$\widehat{\mathbf{w}}^{(t+1)} = \widehat{\mathbf{w}}^{(t)} - \varepsilon \left( \mathbf{E}_{\widehat{\mathbf{w}}^{(t)}}[\mathbf{f} \mid x_{r_t}] - \mathbf{f}(x_{r_t}, y_{r_t}) + 2c\widehat{\mathbf{w}}^{(t)} \right)$$

- Adding  $L_1$  regularisation in SGD and Perceptron introduces *additive weight decay*:

$$\widehat{\mathbf{w}}^{(t+1)} = \widehat{\mathbf{w}}^{(t)} - \varepsilon \left( \mathbf{E}_{\widehat{\mathbf{w}}^{(t)}}[\mathbf{f} \mid x_{r_t}] - \mathbf{f}(x_{r_t}, y_{r_t}) + c \operatorname{sign}(\widehat{\mathbf{w}}^{(t)}) \right)$$



# Stabilising SGD and the Perceptron

- The Perceptron is guaranteed to converge to a weight vector that correctly classifies all training examples *if the training data is separable*
- Most of our problems are *non-separable*
  - ⇒ SGD and the Perceptron never converge to a weight vector
  - ⇒ final weight vector depends on last examples seen
- Reducing learning rate  $\epsilon$  in later iterations can stabilise weight vector  $\hat{\mathbf{w}}$ 
  - ▶ if learning rate is too low, SGD takes a long time to converge
  - ▶ if learning rate is too high,  $\hat{\mathbf{w}}$  can over-shoot
  - ▶ selecting appropriate learning rate is almost “black magic”
- *Bagging* can be used to stabilise SGD and perceptron
  - ▶ *construct multiple models* by running SGD or perceptron many times on random permutations of training data
  - ▶ combine predictions of models at run time by *averaging* or *voting*
- The *averaged perceptron* is a fast approximate version of bagging
  - ▶ train a single perceptron as usual
  - ▶ at end of training, average the weights *from all iterations*
  - ▶ use these averaged weights at run-time

# ADAGRAD and ADADELTA

- There are many methods that attempt to automatically set the learning rate  $\epsilon$
- ADAGRAD and ADADELTA are two of the currently most popular methods
- ADAGRAD estimates a separate learning rate  $\epsilon_j$  for each feature weight  $w_j$
- If  $g_j^{(t)}$  is the derivative of the regularised negative log likelihood  $\ell_D$  w.r.t. feature weight  $w_j$  at step  $t$ , then the ADADGRAD update rule is:

$$\widehat{w}_j^{(t+1)} = \widehat{w}_j^{(t)} - \frac{\eta}{\sqrt{\sum_{t'=1}^t g_j^{(t')}}} g_j^{(t)}, \text{ where:}$$

$$g_j^{(t)} = \frac{\partial \ell_D(\mathbf{w}^{(t)})}{\partial w_j}$$

- This effectively scales the learning rate so features with large derivatives or with fluctuating signs have a slower learning rate
- If a group of features are known to have the same scale, it may make sense for them to share the same learning rate
- The ADADELTA rule is newer and only slightly more complicated
- Both ADAGRAD and ADADELTA only require you to store the sum of the previous derivatives

# Momentum

- Intuition: a ball rolling down a surface will settle at a (local) minimum
- Update should be a mixture of the previous update and derivative of regularised log likelihood  $\ell_D$

$$\begin{aligned}\widehat{\mathbf{w}}^{(t+1)} &= \widehat{\mathbf{w}}^{(t)} + \mathbf{v}^{(t+1)} \\ \mathbf{v}^{(t+1)} &= \alpha \mathbf{v}^{(t)} - (1 - \alpha) \varepsilon \frac{\partial \ell_D(\widehat{\mathbf{w}}^{(t)})}{\partial \mathbf{w}}\end{aligned}$$

- Momentum can smooth statistical fluctuations in SGD derivatives
- If derivatives all point in roughly same direction, updates  $\mathbf{v}$  can become quite large
  - ⇒ set learning rate  $\varepsilon$  to much lower than without momentum
    - ▶ typical value for momentum hyper-parameter  $\alpha = 0.9$



# Comments on SGD and the Perceptron

- Widely used because *easy to implement* and *fast to train*
  - ▶ in my experience, not quite as good as numerical optimisation with L-BFGS
- Overlearning can be a problem
  - ▶ *regularisation* becomes *weight decay*
  - ▶  $L_2$  regularisation is *multiplicative weight decay*
  - ▶  $L_1$  regularisation is *subtractive weight decay*
  - ▶ often more or less ad hoc methods are used instead of or in addition to regularisation
    - averaging (bagging, averaged perceptron, etc.)
    - early stopping
- If you're using either SGD or Perceptron, *try ADAGRAD and ADADELTA learning rules*
  - ▶ these methods *automatically change the learning rate  $\epsilon$*  during learning
  - ▶ they can identify *different learning rates for different features*

⇒ much faster learning than with SGD or Perceptron alone

# Outline

Introducing exponential models

Features in exponential models

Learning exponential models

Regularisation

Conditional models

Stochastic gradient descent and error-driven learning

Avoiding the partition function and its derivatives

Weakly labelled training data

Summary

## Challenges when $\mathcal{Y}$ is large

- The SGD update rule:

$$\widehat{\mathbf{w}}^{(t+1)} = \widehat{\mathbf{w}}^{(t)} - \varepsilon \left( \mathbf{E}_{\widehat{\mathbf{w}}^{(t)}}[\mathbf{f} \mid x_{r_t}] - \mathbf{f}(x_{r_t}, y_{r_t}) \right), \text{ where:}$$

$$\mathbf{E}_{\mathbf{w}}[\mathbf{f} \mid x] = \sum_{y' \in \mathcal{Y}} \mathbf{f}(x, y') \mathbf{P}_{\mathbf{w}}(y' \mid x)$$

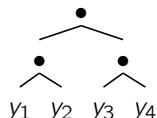
$$\mathbf{P}_{\mathbf{w}}(y \mid x) = \frac{1}{Z(x)} \exp(\mathbf{w} \cdot \mathbf{f}(x, y))$$

$$Z(x) = \sum_{y' \in \mathcal{Y}} \exp(\mathbf{w} \cdot \mathbf{f}(x, y'))$$

- Each update step requires calculating the *partition function*  $Z(x)$  and its derivatives  $\mathbf{E}_{\widehat{\mathbf{w}}^{(t)}}[\mathbf{f} \mid x_{r_t}]$
- These require summing over  $\mathcal{Y}$ , which can dominate the computation time if  $\mathcal{Y}$  is large
  - ▶ in modern speech recognition and machine translation systems,  $\mathcal{Y}$  is the vocabulary of a natural language, so  $|\mathcal{Y}| \approx 10^5$

# Factoring $P(Y | X)$

- Produce a hierarchical clustering of  $\mathcal{Y}$ , which defines a tree over the  $\mathcal{Y}$ .



- Train a separate model for each internal node in the tree
  - ▶ the probability of a leaf (output) node is the *product* of probabilities of each decision on the root to leaf path
- This usually does not produce a very good model
- Conjecture: bagging (e.g., averaging) the output of many such tree models would improve accuracy



# Estimating expected feature counts by sampling

- SGD update rule:

$$\begin{aligned}\widehat{\mathbf{w}}^{(t+1)} &= \widehat{\mathbf{w}}^{(t)} - \varepsilon \left( \mathbf{E}_{\widehat{\mathbf{w}}^{(t)}}[\mathbf{f} \mid x_{r_t}] - \mathbf{f}(x_{r_t}, y_{r_t}) \right), \text{ where:} \\ \mathbf{E}_{\mathbf{w}}[\mathbf{f} \mid x] &= \sum_{y' \in \mathcal{Y}} \mathbf{f}(x, y') \mathbf{P}_{\mathbf{w}}(y' \mid x)\end{aligned}$$

- Idea: use a sampling method to estimate the expected feature counts  $\mathbf{E}_{\widehat{\mathbf{w}}^{(t)}}[\mathbf{f} \mid x_{r_t}]$
- Importance sampling:
  - ▶ draw samples from a *proposal distribution* over  $\mathcal{Y}$  (e.g., unigram distribution)
  - ▶ calculate expectation from samples reweighted according to *importance weights* (which don't require partition function)
- May require a large number of samples to accurately estimate expectations

# Noise-contrastive estimation

- Noise-contrastive estimation can be viewed as *importance sampling with only two samples* (and where importance weights are ignored)
- Suppose the training item at iteration  $t$  is  $(x_{r_t}, y_{r_t}^+)$ .
- Set  $y_t^- \in \mathcal{Y}$  to a random sample from a *proposal distribution* (e.g., unigram distribution over  $\mathcal{Y}$ )
- We approximate:

$$\mathbb{E}_{\widehat{\mathbf{w}}^{(t)}}[\mathbf{f} \mid x_{r_t}] \approx \frac{\mathbf{f}(x_{r_t}, y_{r_t}^+) \exp(\mathbf{w}^{(t)} \cdot \mathbf{f}(x_{r_t}, y_{r_t}^+)) + \mathbf{f}(x_{r_t}, y_t^-) \exp(\mathbf{w}^{(t)} \cdot \mathbf{f}(x_{r_t}, y_t^-))}{\exp(\mathbf{w}^{(t)} \cdot \mathbf{f}(x_{r_t}, y_{r_t}^+)) + \exp(\mathbf{w}^{(t)} \cdot \mathbf{f}(x_{r_t}, y_t^-))}$$

- If  $y_t^-$  is less probable than  $y_{r_t}^+$  the expectation  $\mathbb{E}_{\widehat{\mathbf{w}}^{(t)}}[\mathbf{f} \mid x_{r_t}] \approx \mathbf{f}(x_{r_t}, y_{r_t}^+)$ , so the expectations will cancel, and there won't be a large weight update
- If  $y_t^-$  is more probable than  $y_{r_t}^+$  the expectation  $\mathbb{E}_{\widehat{\mathbf{w}}^{(t)}}[\mathbf{f} \mid x_{r_t}] \approx \mathbf{f}(x_{r_t}, y_t^-)$ , so there can be a large weight update
- Widely used in the neural net community today

# Outline

Introducing exponential models

Features in exponential models

Learning exponential models

Regularisation

Conditional models

Stochastic gradient descent and error-driven learning

Avoiding the partition function and its derivatives

**Weakly labelled training data**

Summary

## Ambiguous or weakly labelled training data as partial observations

- Suppose our training data doesn't tell us the *true label*  $y_i$  for each example  $x_i$ , but only provides us with *a set of labels*  $\mathcal{Y}_i$  that contains the unknown true label  $y_i$

$$D = ((x_1, \mathcal{Y}_1), \dots, (x_n, \mathcal{Y}_n)) \text{ where:}$$
$$y_i \in \mathcal{Y}_i \subseteq \mathcal{Y}$$

- Idea: learn a model that maximizes  $\prod_{i=1}^n \mathbf{P}(\mathcal{Y}_i | x_i)$
- Example: in reranking the gold parse might not be in the beam, so train model to select one of the best parses available in beam; we don't care which is chosen
- Example: in arc-eager dependency parsing, several different moves can lead to same gold parse; we don't care which the parser chooses

## Partially-observed conditional exponential models

- Data  $D = (((x_1, \mathcal{Y}_1), \dots, (x_n, \mathcal{Y}_n)))$ , where  $\mathcal{Y}_i \subseteq \mathcal{Y}$  and  $\mathcal{Y}_i \neq \emptyset$
- Compute *maximum conditional likelihood estimator* by minimizing negative log conditional likelihood

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \ell_D(\mathbf{w}), \text{ where:}$$

$$\ell_D(\mathbf{w}) = - \sum_{i=1}^n \log \mathbf{P}_{\mathbf{w}}(\mathcal{Y}_i | x_i)$$

$$= \sum_{i=1}^n (\log Z(x_i, \mathcal{Y}) - \log Z(x_i, \mathcal{Y}_i)), \text{ where:}$$

$$Z(x, \mathcal{Y}') = \sum_{y' \in \mathcal{Y}'} \exp(\mathbf{w} \cdot \mathbf{f}(x, y'))$$

- Intuition:  $\log Z(x_i, \mathcal{Y}) - \log Z(x_i, \mathcal{Y}_i)$  will be small when most mass is assigned to  $\mathcal{Y}_i$
- If  $\mathcal{Y}_i = \mathcal{Y}$ , then example  $i$  *has no information*
- Warning:  $\ell_D$  is usually *not convex*  $\Rightarrow$  local minima
  - ▶ hidden data problems usually have non-convex log likelihoods

# Derivatives for partially-observed conditional models

- Negative log likelihood:

$$\ell_D(\mathbf{w}) = \sum_{i=1}^n (\log Z(x_i, \mathcal{Y}) - \log Z(x_i, \mathcal{Y}_i)), \text{ where:}$$

$$Z(x, \mathcal{Y}') = \sum_{y' \in \mathcal{Y}'} \exp(\mathbf{w} \cdot \mathbf{f}(x, y'))$$

- Derivatives are differences of *two conditional expectations*

$$\frac{\partial \ell_D}{\partial w_j} = \sum_{i=1}^n (\mathbf{E}_{\mathbf{w}}[f_j | x_i, \mathcal{Y}] - \mathbf{E}_{\mathbf{w}}[f_j | x_i, \mathcal{Y}_i]), \text{ where:}$$

$$\mathbf{E}_{\mathbf{w}}[f_j | x, \mathcal{Y}'] = \sum_{y' \in \mathcal{Y}'} f_j(x, y') \mathbf{P}_{\mathbf{w}}(y' | x) \quad (\text{expected value given } x \text{ and } \mathcal{Y}')$$

- These derivatives are no harder to compute than for the fully-observed case
- SGD and perceptron algorithms generalise straight-forwardly to partially-observed data

# Partially-observed data in the reranker

- Training data consists of a sequence of *training data items* (sentences)
- Each data item consists of a sequence of *candidates* (parses)
  - ▶ the number of candidates per data item can vary
- Each candidate consists of a sequence of *feature-value pairs*
- Each feature is an integer, and each value is a floating-point number
  - ▶ feature value 1 is special-cased because it's so common in “one-hot” representations
- To allow partially-observed training data, each candidate has a *gold weight*
  - ▶ for a standard MaxEnt model, the gold candidate in each data item has gold weight 1, all others have gold weight 0
  - ▶ with partially-observed data, more than one candidate has weight 1

## Other interesting things the reranker can do

- Data items (sentences) and candidates (parses) can be given “costs” so the reranker can calculate f-scores
  - ▶ can optimise *expected f-score* instead of log likelihood
  - ▶ useful with skewed data (e.g., in disfluency detection, where most words are fluent)
- The reranker uses L1 and/or L2 regularisation
  - ▶ can optimise regulariser constants to maximise log likelihood or f-score of heldout data
- Features are organised into *feature classes*
  - ▶ each feature class can have its own regulariser constant
  - ▶ these feature constants can be optimised can be on heldout data
- Standard optimiser is L-BFGS-OWLQN, but can also use Averaged Perceptron
  - ▶ Averaged Perceptron is not quite as good as L-BFGS, but *much faster*
  - ▶ Averaged Perceptron can be used to search for subset of feature classes that optimise f-score on heldout data



# Outline

Introducing exponential models

Features in exponential models

Learning exponential models

Regularisation

Conditional models

Stochastic gradient descent and error-driven learning

Avoiding the partition function and its derivatives

Weakly labelled training data

Summary

# Summary

- Maximum entropy models capture the intuition that *features interact multiplicatively*, i.e., can increase or decrease the probability of an outcome
- Calculating the *partition function*  $Z$  and its derivatives is usually the central challenge in MaxEnt modelling
- *Conditional* MaxEnt models, which model  $\mathbf{P}(y \mid \mathbf{x})$ , often have simpler partition functions than *joint* models, which model  $\mathbf{P}(y, \mathbf{x})$ .
- *Regularisation* is often essential to avoid over-learning
  - ▶ *L1 regularisation* produces sparse feature weight vectors  $\mathbf{w}$  at the individual feature level
  - ▶ the *group Lasso* produces sparse feature weight vectors  $\mathbf{w}$  at the feature group level
- Stochastic Gradient Descent (SGD) is an easy and fast way to learn MaxEnt models (but less accurate?)
- The Perceptron is SGD for conditional MaxEnt with a Viterbi approximation

# Where we go from here

- *Conditional Random Fields* are conditional MaxEnt models that use dynamic programming to calculate the partition function and their derivatives
  - ▶ generally requires  $\mathcal{Y}$  to have some kind of regular structure, e.g., a sequence (sequence labelling) or a tree (parsing)
- *Neural networks* use MaxEnt models as components (they are networks of MaxEnt models, but “neural net” sounds better!)
  - ▶ *Boltzmann machines* are MaxEnt models where the data items are graphs
  - ▶ *feed-forward networks* use conditional MaxEnt models as components