Lecture 3: Multi-Class Classification

Kai-Wei Chang CS @ UCLA <u>kw@kwchang.net</u>

Couse webpage: https://uclanlp.github.io/CS269-17/



ML in NLP

1

Previous Lecture

Binary linear classification models

- Perceptron, SVMs, Logistic regression
- Prediction is simple:
 - Given an example x, prediction is $sgn(w^Tx)$
 - Note that all these linear classifier have the same inference rule
 - In logistic regression, we can further estimate the probability

$$P(y=1|\mathbf{x},\mathbf{w}) = \frac{e^{\mathbf{w}^T \mathbf{x}}}{1+e^{\mathbf{w}^T \mathbf{x}}} = \frac{1}{1+e^{-\mathbf{w}^T \mathbf{x}}}$$

Question?

This Lecture

Multiclass classification overview

- Reducing multiclass to binary
 - One-against-all & One-vs-one
 - Error correcting codes
- Training a single classifier
 - Multiclass Perceptron: Kesler's construction
 - Multiclass SVMs: Crammer&Singer formulation
 - Multinomial logistic regression



What is multiclass

♦ Output \in {1,2,3, ... *K*}

In some cases, output space can be very large (i.e., K is very large)

Each input belongs to exactly one class
 (c.f. in multilabel, input belongs to many classes)





Multi-class Applications in NLP?



Two key ideas to solve multiclass

Reducing multiclass to binary

- Decompose the multiclass prediction into multiple binary decisions
- Make final decision based on multiple binary classifiers
- Training a single classifier
 - Minimize the empirical risk
 - Consider all classes simultaneously



Reduction v.s. single classifier

Reduction

Future-proof: binary classification improved so does muti-class

Easy to implement

Single classifier

- Global optimization: directly minimize the empirical loss; easier for joint prediction
- Easy to add constraints and domain knowledge



A General Formula

 $\hat{y} = \operatorname{argmax}_{y \in \mathcal{Y}} f(y; w, x)$ input model parameters output space Inference/Test: given w, x, solve argmax Learning/Training: find a good w ♦ Today: $x \in \mathbb{R}^n$, $\mathcal{Y} = \{1, 2, ..., K\}$ (multiclass)



This Lecture

Multiclass classification overview

- Reducing multiclass to binary
 - One-against-all & One-vs-one
 - Error correcting codes
- Training a single classifier
 - Multiclass Perceptron: Kesler's construction
 - Multiclass SVMs: Crammer&Singer formulation
 - Multinomial logistic regression



One against all strategy





CS6501 Lecture 3

One against All learning

✤ Multiclass classifier
♦ Function f: Rⁿ → {1,2,3,...,k}

Decompose into binary problems







One-again-All learning algorithm

- ✤ Learning: Given a dataset $D = \{(x_i, y_i)\}$ $x_i \in \mathbb{R}^n, y_i \in \{1, 2, 3, ..., K\}$
- Decompose into K binary classification tasks
 - ***** Learn K models: $w_1, w_2, w_3, \dots w_K$
 - For class k, construct a binary classification task as:
 - Positive examples: Elements of D with label k
 - Negative examples: All other elements of D
 - The binary classification can be solved by any algorithm we have seen



One against All learning

Multiclass classifier

♦ Function $f : \mathbb{R}^n \rightarrow \{1, 2, 3, ..., k\}$



Decompose into binary problems

Ideal case: only the correct label will have a positive score





CS6501 Lecture 3

One-again-All Inference

★ Learning: Given a dataset $D = \{(x_i, y_i)\}$ $x_i \in \mathbb{R}^n, y_i \in \{1, 2, 3, ..., K\}$

Decompose into K binary classification tasks
 Learn K models: w₁, w₂, w₃, ... w_K

Inference: "Winner takes all"

$$\widehat{y} = \operatorname{argmax}_{y \in \{1, 2, \dots, K\}} W_y^T x$$

For example: $y = \operatorname{argmax}(w_{black}^T x, w_{blue}^T x, w_{green}^T x)$

An instance of the general form

$$\hat{y} = \operatorname{argmax}_{y \in \mathcal{Y}} f(y; w, x)$$

$$w = \{w_1, w_2, ..., w_K\}, f(y; w, x) = w_y^T x$$



One-again-All analysis

Not always possible to learn

- Assumption: each class individually separable from all the others
- No theoretical justification
 - Need to make sure the range of all classifiers is the same – we are comparing scores produced by K classifiers trained independently.
- Easy to implement; work well in practice





One v.s. One (All against All) strategy





CS6501 Lecture 3

One v.s. One learning

♦ Multiclass classifier
♦ Function f: Rⁿ → {1,2,3,...,k}

Decompose into binary problems







One-v.s-One learning algorithm

- ✤ Learning: Given a dataset $D = \{(x_i, y_i)\}$ $x_i \in \mathbb{R}^n, y_i \in \{1, 2, 3, ..., K\}$
- Decompose into C(K,2) binary classification tasks
 - * Learn C(K,2) models: $w_1, w_2, w_3, ..., w_{K^*(K-1)/2}$
 - For each class pair (i,j), construct a binary classification task as:
 - Positive examples: Elements of D with label i
 - Negative examples Elements of D with label j
 - The binary classification can be solved by any algorithm we have seen



One-v.s-One Inference algorithm

- Decision Options:
 - More complex; each label gets k-1 votes
 - Output of binary classifier may not cohere.
 - Majority: classify example x to take label i if i wins on x more often than j (j=1,...k)
 - A tournament: start with n/2 pairs; continue with winners



Classifying with One-vs-one



All are post-learning and *might* cause weird stuff



CS6501 Lecture 3

One-v.s.-one Assumption

Every pair of classes is separable





CS6501 Lecture 3

Comparisons

- One against all
 - O(K) weight vectors to train and store
 - Training set of the binary classifiers may unbalanced
 - Less expressive; make a strong assumption
- One v.s. One (All v.s. All)
 - $O(K^2)$ weight vectors to train and store
 - Size of training set for a pair of labels could be small
 ⇒ overfitting of the binary classifiers
 - Need large space to store model



Problems with Decompositions

- Learning optimizes over *local* metrics
 - Does not guarantee good global performance
 - We don't care about the performance of the local classifiers
- Poor decomposition ⇒ poor performance
 - Difficult local problems
 - Irrelevant local problems
- Efficiency: e.g., All vs. All vs. One vs. All
- Not clear how to generalize multi-class to problems with a very large # of output



Still an ongoing research direction

Key questions:

- How to deal with large number of classes
- How to select "right samples" to train binary classifiers

Error-correcting tournaments

[Beygelzimer, Langford, Ravikumar 09]

Logarithmic Time One-Against-Some

[Daume, Karampatziakis, Langford, Mineiro 16]

Label embedding trees for large multi-class tasks. [Bengio, Weston, Grangier 10]



Decomposition methods: Summary

General Ideas:

- Decompose the multiclass problem into many binary problems
- Prediction depends on the decomposition
 - Constructs the multiclass label from the output of the binary classifiers
- Learning optimizes local correctness
 - Each binary classifier don't need to be globally correct and isn't aware of the prediction procedure



This Lecture

Multiclass classification overview

- Reducing multiclass to binary
 - One-against-all & One-vs-one
 - Error correcting codes
- Training a single classifier
 - Multiclass Perceptron: Kesler's construction
 - Multiclass SVMs: Crammer&Singer formulation
 - Multinomial logistic regression



Revisit One-again-All learning algorithm

- ✤ Learning: Given a dataset $D = \{(x_i, y_i)\}$ $x_i \in \mathbb{R}^n, y_i \in \{1, 2, 3, ..., K\}$
- Decompose into K binary classification tasks
 - ***** Learn K models: $w_1, w_2, w_3, \dots w_K$
 - v_k : separate class k from others

Prediction

$$\hat{y} = \operatorname{argmax}_{y \in \{1,2,\ldots K\}} w_y^T x$$



Observation

At training time, we require $w_i^T x$ to be positive for examples of class *i*.

✤ Really, all we need is for $w_i^T x$ to be more than all others ⇒ this is a weaker requirement

For examples with label *i*, we need $w_i^T x > w_j^T x$ for all *j*



Perceptron-style algorithm

For examples with label *i*, we need $w_i^T x > w_j^T x$ for all *j*

• For each training example (x, y)

- ♦ If for some y', $w_y^T x \le w_{y'}^T x$ mistake!
 - $w_y \leftarrow w_y + \eta x$ update to promote y

 $\bigstar w_{y'} \leftarrow w_{y'} - \eta x$

update to demote y'

Why add ηx to w_y promote label y: Before update $s(y) = \langle w_y^{old}, x \rangle$ After update $s(y) = \langle w_y^{new}, x \rangle = \langle w_y^{old} + \eta x, x \rangle$ $= \langle w_y^{old}, x \rangle + \eta \langle x, x \rangle$ Note! $\langle x, x \rangle = x^T x > 0$



A Perceptron-style Algorithm

Given a training set $\mathcal{D} = \{(x, y)\}$ Initialize $w \leftarrow \mathbf{0} \in \mathbb{R}^n$ For epoch 1...T: How to analyze this algorithm For (x, y) in \mathcal{D} : and simplify the update rules? For $y' \neq y$ if $w_v^T x < w_{v'}^T x$ make a mistake $w_{v} \leftarrow w_{v} + \eta x$ promote y demote y' $w_{\gamma\prime} \leftarrow w_{\gamma\prime} - \eta x$ Return w

Prediction: $\operatorname{argmax}_{y} w_{y}^{T} x$



Linear Separability with multiple classes

Let's rewrite the equation

 $w_i^T x > w_j^T x$ for all j

✤ Instead of having $w_1, w_2, w_3, ..., w_K$, we want to represent the model using a single vector w

 w^T ? > w^T ? for all j

✤ How?

Change the input representation Let's define $\phi(x, y)$, such that $w^T \phi(x, i) > w^T \phi(x, j) \quad \forall j$

multiple models v.s. multiple data points

Assume we have a multi-class problem with K class and n features.

$$w_i^T x > w_j^T x \quad \forall j$$

$$\bigstar \text{ models:}$$

$$w_1, w_2, \dots w_K, \qquad w_k \in R^n$$

$$\bigstar \text{ Input:}$$

$$x \in R^n$$

 $w^T \phi(x,i) > w^T \phi(x,j) \quad \forall j$



Assume we have a multi-class problem with K class and n features.

$$w_i^T x > w_j^T x \quad \forall j$$

$$\bigstar \text{ models:}$$

$$w_1, w_2, \dots w_K, \qquad w_k \in \mathbb{R}^n$$

$$\diamondsuit \text{ Input:}$$

$$x \in \mathbb{R}^n$$

 $w^{T} \phi(x, i) > w^{T} \phi(x, j) \quad \forall j$ Only one model: $w \in R^{n \times K}$





Assume we have a multi-class problem with K class and n features.

$$w_i^T x > w_j^T x \quad \forall j$$

$$\bigstar \text{ models:}$$

$$w_1, w_2, \dots w_K, \qquad w_k \in \mathbb{R}^n$$

$$\bigstar \text{ lnput:}$$

$$x \in \mathbb{R}^n$$

$$w^T \phi(x,i) > w^T \phi(x,j) \quad \forall j$$

- Only one model: $w \in R^{nK}$
- Define $\phi(x, y)$ for label y being associated to input x

$$(x, y) = \begin{bmatrix} 0_n \\ \vdots \\ x \\ \vdots \\ 0_n \end{bmatrix}_{nK \times 1} x \text{ in } y^{th} \text{ block;}$$

Zeros elsewhere



φ

Assume we have a multi-class problem with K class and n features.

$$w^{T} \phi(x, i) > w^{T} \phi(x, j) \quad \forall j$$

$$w = \begin{bmatrix} W_{1} \\ \vdots \\ W_{y} \\ \vdots \\ W_{n} \end{bmatrix}_{nK \times 1} \qquad \phi(x, y) = \begin{bmatrix} 0_{n} \\ \vdots \\ x \\ \vdots \\ 0_{n} \end{bmatrix}_{nK \times 1}$$

$$x \text{ in } y^{th} \text{ block};$$
Zeros elsewhere
$$x \\ \vdots \\ 0_{n} \end{bmatrix}_{nK \times 1}$$

$$w^T \phi(x, y) = w_y^T x$$



CS6501 Lecture 3

Assume we have a multi-class problem with K class and n features.




Linear Separability with multiple classes



 $\phi(\mathbf{x}, i) = \begin{bmatrix} n \\ \vdots \\ \mathbf{x} \leftarrow - \mathbf{j}^{th} \\ \vdots \\ \mathbf{0}_n \end{bmatrix}_{nK \times 1}^{n}$ block

Positive examples

$$\phi(\mathbf{x},i) - \phi(\mathbf{x},j)$$

Negative examples

 $-\phi(\mathbf{x},i) + \phi(\mathbf{x},j)$



CS6501 Lecture 3

How can we predict?

$$\operatorname{argmax}_{y} w^{T} \phi(x, y)$$

$$w = \begin{bmatrix} w_1 \\ \vdots \\ w_y \\ \vdots \\ w_n \end{bmatrix}_{nK \times 1} \phi(x, y) = \begin{bmatrix} 0_n \\ \vdots \\ x \\ \vdots \\ 0_n \end{bmatrix}_{nK \times 1}$$

For input an input x, the model predict label is 3





How can we predict?

$$\operatorname{argmax}_{y} w^{T} \phi(x, y)$$



For input an input x, the model predict label is 3



This is equivalent to $\operatorname{argmax}_{y \in \{1,2,\dots K\}} W_y^T x$



Constraint Classification

• Goal: $w[\phi(x,i) - \phi(x,j)] \ge 0 \quad \forall j$

Training:

• For each example (x, i)

♦ Update model if $w^T[\phi(x,i) - \phi(x,j)] < 0$, $\forall j$



Given a training set $\mathcal{D} = \{(x, y)\}$ Initialize $w \leftarrow \mathbf{0} \in \mathbb{R}^n$ For epoch $1 \dots T$: This needs $|D| \times K$ updates, For (x, y) in \mathcal{D} : do we need all of them? For $y' \neq y$ if $w^{T}[\phi(x, y) - \phi(x, y')] < 0$ $w \leftarrow w + \eta \left[\phi(x, y) - \phi(x, y')\right]$ Return w How to interpret this update rule? $\operatorname{argmax}_{v} w^{T} \phi(x, y)$ Prediction:



An alternative training algorithm



***** For each example (x, i)

Find the prediction of the current model:

 $\hat{y} = \operatorname{argmax}_{j} w^{T} \phi(x, j)$

♦ Update model if $w^T[\phi(x,i) - \phi(x,\hat{y})] < 0$, $\forall y'$



Given a training set $\mathcal{D} = \{(x, y)\}$ Initialize $w \leftarrow \mathbf{0} \in \mathbb{R}^n$ For epoch 1...T: For (x, y) in \mathcal{D} : $\hat{y} = \operatorname{argmax}_{v'} w^T \phi(x, y')$ if $w^{T}[\phi(x, y) - \phi(x, \hat{y})] < 0$ $w \leftarrow w + \eta \left[\phi(x, y) - \phi(x, \widehat{y}) \right]$ Return w How to interpret this update rule? **Prediction:** $\operatorname{argmax}_{v} w^{T} \phi(x, y)$



Given a training set $\mathcal{D} = \{(x, y)\}$ Initialize $w \leftarrow \mathbf{0} \in \mathbb{R}^n$ There are only two situations: For epoch 1...T: 1. $\hat{y} = y$: $\phi(x, y) - \phi(x, \hat{y}) = 0$ 2. $w^T[\phi(x, y) - \phi(x, \hat{y})] < 0$ For (x, y) in \mathcal{D} : $\hat{y} = \operatorname{argmax}_{v}, w^T \phi(x, y')$ if $w^T[\phi(x,y) - \phi(x,\hat{y})] < 0$ $w \leftarrow w + \eta \left[\phi(x, y) - \phi(x, \hat{y}) \right]$ Return w How to interpret this update rule? $\operatorname{argmax}_{v} w^{T} \phi(x, y)$ Prediction:



Given a training set $\mathcal{D} = \{(x, y)\}$ Initialize $w \leftarrow \mathbf{0} \in \mathbb{R}^n$ For epoch $1 \dots T$: For (x, y) in \mathcal{D} : $\hat{y} = \operatorname{argmax}_{y}, w^T \phi(x, y')$ $w \leftarrow w + \eta [\phi(x, y) - \phi(x, \hat{y})]$ Return w

How to interpret this update rule?

Prediction: $\operatorname{argmax}_{y} w^{T} \phi(x, y)$



Consider multiclass margin





Marginal constraint classifier

Goal: for every (x,y) in the training data set

 $\min_{y' \neq y} w^{T} [\phi(x, y) - \phi(x, y')] \ge \delta$ $\Rightarrow w^{T} \phi(x, y) - \max_{y \neq y'} w^{T} \phi(x, y') \ge \delta$ $\Rightarrow w^{T} \phi(x, i) - [\max_{y \neq y'} w^{T} \phi(x, j) + \delta] \ge 0$

Score for a label Labels

 ${f A}$ engineering

Computer Science

Let's define: $\Delta(y, y') = \begin{cases} \delta & \text{if } y \neq y' \\ 0 & if \ y = y' \end{cases}$ Check if $y = \arg \max_{y'} w^T \phi(x, y') + \Delta(y, y')$

Constraints violated \Rightarrow need an update

CS6501 Lecture 3

Given a training set $\mathcal{D} = \{(x, y)\}$ Initialize $w \leftarrow \mathbf{0} \in \mathbb{R}^n$ **Multiclass Margin** δ Blue Score for For epoch 1...T: Red a label Green Black For (x, y) in \mathcal{D} : Labels $\hat{y} = \operatorname{argmax}_{y'} w^T \phi(x, y') + \Delta(y, y')$ $w \leftarrow w + \eta \left[\phi(x, y) - \phi(x, \hat{y}) \right]$ Return w

How to interpret this update rule?

Prediction: $\operatorname{argmax}_{y} w^{T} \phi(x, y)$



Remarks

This approach can be generalized to train a ranker; in fact, any output structure

We have preference over label assignments

E.g., rank search results, rank movies / products



A peek of a generalized Perceptron model

Given a training set $\mathcal{D} = \{(x, y)\}$ Initialize $w \leftarrow \mathbf{0} \in \mathbb{R}^n$ Structured output For epoch $1 \dots T$: For (x, y) in \mathcal{D} : Structural prediction/Inference $\hat{y} = \operatorname{argmax}_{y'} w^T \phi(x, y') + \Delta(y, y')$ $\boldsymbol{w} \leftarrow \boldsymbol{w} + \eta \left[\phi(x, y) - \phi(x, \hat{y}) \right]$ Structural loss Return w Model update $\operatorname{argmax}_{v} w^{T} \phi(x, y)$ Prediction:



Recap: A Perceptron-style Algorithm

Given a training set $\mathcal{D} = \{(x, y)\}$ Initialize $w \leftarrow \mathbf{0} \in \mathbb{R}^n$ For epoch $1 \dots T$: For (x, y) in \mathcal{D} : For $y' \neq y$ if $w_v^T x < w_{v'}^T x$ make a mistake $w_{v} \leftarrow w_{v} + \eta x$ promote y demote y' $w_{y'} \leftarrow w_{y'} - \eta x$ Return w

Prediction: $\operatorname{argmax}_{y} w_{y}^{T} x$



Recap: Kesler construction

Assume we have a multi-class problem with K class and n features.

$$w^{T} \phi(x, i) > w^{T} \phi(x, j) \quad \forall j$$

$$w = \begin{bmatrix} W_{1} \\ \vdots \\ W_{y} \\ \vdots \\ W_{n} \end{bmatrix}_{nK \times 1} \qquad \phi(x, y) = \begin{bmatrix} 0_{n} \\ \vdots \\ x \\ \vdots \\ 0_{n} \end{bmatrix}_{nK \times 1}$$

$$x \text{ in } y^{th} \text{ block; } Zeros \text{ elsewhere}$$

$$w^T \phi(x, y) = w_y^T x$$



CS6501 Lecture 3





CS6501 Lecture 3

Recap: A Perceptron-style Algorithm

Given a training set $\mathcal{D} = \{(x, y)\}$ Initialize $w \leftarrow \mathbf{0} \in \mathbb{R}^n$ For epoch $1 \dots T$: For (x, y) in \mathcal{D} : $\hat{y} = \operatorname{argmax}_{y}, w^T \phi(x, y')$ $w \leftarrow w + \eta [\phi(x, y) - \phi(x, \hat{y})]$ Return w

How to interpret this update rule?

Prediction: $\operatorname{argmax}_{y} w^{T} \phi(x, y)$



Multi-category to Constraint Classification

Multiclass

 $(\mathsf{x}, \mathsf{A}) \implies (\mathsf{x}, (\mathsf{A} \ge \mathsf{B}, \mathsf{A} \ge \mathsf{C}, \mathsf{A} \ge \mathsf{D}))$

Multilabel

 $\diamondsuit (\mathsf{x}, (\mathsf{A}, \mathsf{B})) \Rightarrow (\mathsf{x}, ((\mathsf{A} \ge \mathsf{C}, \mathsf{A} \ge \mathsf{D}, \mathsf{B} \ge \mathsf{C}, \mathsf{B} \ge \mathsf{D}))$

Label Ranking

 $(x, (5>4>3>2>1)) \implies (x, ((5>4, 4>3, 3>2, 2>1))$



Generalized constraint classifiers

❖ In all cases, we have examples (x,y) with $y \in S_k$

- Where S_k : partial order over class labels {1,...,k}
 defines "preference" relation (>) for class labeling
- ✤ Consequently, the Constraint Classifier is: h: $X \rightarrow S_k$ ♠ h(x) is a partial order

♦ h(x) is *consistent* with y if $(i < j) \in y \rightarrow (i < j) \in h(x)$

Just like in the multiclass we learn one $w_i \in R^n$ for each label, the same is done for multi-label and ranking. The weight vectors are updated according with the requirements from $y \in S_k$



Multi-category to Constraint Classification

Solving structured prediction problems by ranking algorithms

Multiclass

 $\bigstar (\mathsf{x}, \mathsf{A}) \implies (\mathsf{x}, (\mathsf{A} \ge \mathsf{B}, \mathsf{A} \ge \mathsf{C}, \mathsf{A} \ge \mathsf{D}))$

Multilabel

 $(\mathsf{x}, (\mathsf{A}, \mathsf{B})) \Rightarrow (\mathsf{x}, ((\mathsf{A} \ge \mathsf{C}, \mathsf{A} \ge \mathsf{D}, \mathsf{B} \ge \mathsf{C}, \mathsf{B} \ge \mathsf{D}))$

Label Ranking

 $\begin{array}{ll} \bigstar \ (x, \ (5 > 4 > 3 > 2 > 1)) & \Rightarrow (x, \ (\ (5 > 4, \ 4 > 3, \ 3 > 2, \ 2 > 1)) \\ y \in {\textbf{S}}_k & h: {\textbf{X}} \to {\textbf{S}}_k \end{array}$



Properties of Construction (Zimak et. al 2002, 2003)

- Can learn any argmax w_i.x function (even when i isn't linearly separable from the union of the others)
- Can use any algorithm to find linear separation
 - Perceptron Algorithm
 - *ultraconservative online algorithm* [Crammer, Singer 2001]
 - Winnow Algorithm
 - *multiclass winnow* [Masterharm 2000]
- Defines a *multiclass margin* by binary margin in R^{KN}
 multiclass SVM [Crammer, Singer 2001]



This Lecture

Multiclass classification overview

- Reducing multiclass to binary
 - One-against-all & One-vs-one
 - Error correcting codes
- Training a single classifier
 - Multiclass Perceptron: Kesler's construction
 - Multiclass SVMs: Crammer&Singer formulation
 - Multinomial logistic regression



Recall: Margin for binary classifiers

The margin of a hyperplane for a dataset is the distance between the hyperplane and the data point nearest to it.



Multi-class SVM

In a risk minimization framework

• Goal: $D = \{(x_i, y_i)\}_{i=1}^N$

- **1.** $w_{y_i}^T x_i > w_{y'}^T x_i$ for all *i*, y'
- 2. Maximizing the margin

Defined as the score difference between the highest scoring label and the second one





Multiclass Margin

Defined as the score difference between the highest scoring label and the second one



Multiclass SVM (Intuition)

Binary SVM

Maximize margin. Equivalently,

Minimize norm of weights such that the closest points to the hyperplane have a score 1

Multiclass SVM

- Each label has a different weight vector (like one-vs-all)
- Maximize multiclass margin. Equivalently,

Minimize total norm of the weights such that the true label is scored at least 1 more than the second best one

Multiclass SVM in the separable case

Computer Science

Multiclass SVM: General case

Multiclass SVM: General case

Recap: An alternative SVM formulation

$$\min_{\mathbf{w},b,\xi} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \xi_i$$

s.t $y_i (\mathbf{w}^T \mathbf{x}_i + b) \ge 1 - \xi_i; \ \xi_i \ge 0 \quad \forall i$
* Rewrite the constraints:
 $\xi_i \ge 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b); \ \xi_i \ge 0 \quad \forall i$
* In the optimum, $\xi_i = \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b))$
* Soft SVM can be rewritten as:
$$\min_{\mathbf{w},b} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + \mathbf{w}_i + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + \mathbf{w}_i + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + \mathbf{w}_i + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + \mathbf{w}_i + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + \mathbf{w}_i + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + \mathbf{w}_i + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + \mathbf{w}_i + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + \mathbf{w}_i + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + \mathbf{w}_i + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + \mathbf{w}_i + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + \mathbf{w}_i + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + \mathbf{w}_i + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + \mathbf{w}_i + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + \mathbf{w}_i + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + \mathbf{w}_i + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + \mathbf{w}_i + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + \mathbf{w}_i + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + \mathbf{w}_i + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + \mathbf{w}_i + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + \mathbf{w}_i + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + \mathbf{w}_i + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + \mathbf{w}_i + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + \mathbf{w}_i + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + \mathbf{w}_i + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + \mathbf{w}_i + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + \mathbf{w}_i + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b) + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b) + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b) + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b) + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b) + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b) + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b) + C \sum_i \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)$$

Rewrite it as unconstraint problem

$$\min_{\mathbf{w}_1, \mathbf{w}_2, \cdots, \mathbf{w}_K, \xi} \quad \frac{1}{2} \sum_k \mathbf{w}_k^T \mathbf{w}_k + C \sum_{(\mathbf{x}_i, \mathbf{y}_i) \in D} \xi_i$$
s.t.
$$\mathbf{w}_{\mathbf{y}_i}^T \mathbf{x} - \mathbf{w}_k^T \mathbf{x} \ge 1 - \xi_i, \qquad \forall (\mathbf{x}_i, \mathbf{y}_i) \in D,$$

$$k \in \{1, 2, \cdots, K\}, k \neq \mathbf{y}_i,$$

$$\xi_i \ge 0, \qquad \forall i.$$

Let's define: $\Delta(y, y') = \begin{cases} \delta & \text{if } y \neq y' \\ 0 & if \ y = y' \end{cases}$

 $\min_{w} \frac{1}{2} \sum_{k} w_{k}^{T} w_{k} + C \sum_{i} (\max_{k} (\Delta(y_{i}, k) + w_{k}^{T} x) - w_{y_{i}}^{T} x)$

Multiclass SVM

Generalizes binary SVM algorithm

- If we have only two classes, this reduces to the binary (up to scale)
- Comes with similar generalization guarantees as the binary SVM
- Can be trained using different optimization methods
 - Stochastic sub-gradient descent can be generalized

Write down SGD for multiclass SVM

Write down multiclas SVM with Kesler construction

This Lecture

Multiclass classification overview

- Reducing multiclass to binary
 - One-against-all & One-vs-one
 - Error correcting codes
- Training a single classifier
 - Multiclass Perceptron: Kesler's construction
 - Multiclass SVMs: Crammer&Singer formulation
 - Multinomial logistic regression

Recall: (binary) logistic regression

$$\min_{\boldsymbol{w}} \quad \frac{1}{2}\boldsymbol{w}^{T}\boldsymbol{w} + C\sum_{i} \log(1 + e^{-y_{i}(\boldsymbol{w}^{T}\boldsymbol{x}_{i})})$$

Assume labels are generated using the following probability distribution:

$$P(y = 1 | \mathbf{x}, \mathbf{w}) = \frac{e^{\mathbf{w}^T \mathbf{x}}}{1 + e^{\mathbf{w}^T \mathbf{x}}} = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$$
$$P(y = -1 | \mathbf{x}, \mathbf{w}) = \frac{1}{1 + e^{\mathbf{w}^T \mathbf{x}}}$$

(multi-class) log-linear model



- This is a valid probability assumption. Why?
- Another way to write this (with Kesler construction) is
 This often

This often called soft-max

$$P(y|x,w) = \frac{\exp(w^T \phi(x,y))}{\sum_{y' \in \{1,2,...K\}} \exp(w^T \phi(x,y'))}$$



Softmax

Softmax: let s(y) be the score for output y here s(y)= $w^T \phi(x, y)$ (or $w_y^T x$) but it can be computed by other metric. $P(y) = \frac{\exp(s(y))}{\sum_{y' \in \{1, 2, \dots, K\}} \exp(s(y))}$

We can control the peakedness of the distribution

$$P(y|\sigma) = \frac{\exp(s(y)/\sigma)}{\sum_{y' \in \{1,2,\dots,K\}} \exp(s(y/\sigma))}$$



Example



Log linear model



Note:

 $p(y) \propto \exp(\theta^{\mathsf{T}} \mathbf{f}(y))$



Maximum log-likelihood estimation

Training can be done by maximum log-likelihood estimation i.e. $\max_{w} \log P(D|w)$

 $D = \{(x_i, y_i)\}$

$$P(D|w) = \prod_{i} \frac{\exp(w_{y_{i}}^{T} x_{i})}{\sum_{y' \in \{1,2,\dots,K\}} \exp(w_{y'}^{T} x_{i})}$$
$$\log P(D|w) = \sum_{i} [w_{y_{i}}^{T} x_{i} - \log \sum_{y' \in \{1,2,\dots,K\}} \exp(w_{y'}^{T} x_{i})]$$



Maximum a posteriori

 $D = \{(x_i, y_i)\}$

Can you use Kesler construction to rewrite this formulation?

 $P(w|D) \propto P(w)P(D|w)$

$$\max_{w} -\frac{1}{2} \sum_{y} w_{y}^{T} w_{y} + C \sum_{i} [w_{y_{i}}^{T} x_{i} - \log \sum_{y' \in \{1, 2, \dots, K\}} \exp(w_{y'}^{T} x_{i})]$$

or

$$\min_{w} \frac{1}{2} \sum_{y} w_{y}^{T} w_{y} + C \sum_{i} [\log \sum_{y' \in \{1,2,\dots,K\}} \exp(w_{y'}^{T} x_{i}) - w_{y_{i}}^{T} x_{i}]$$



Comparisons

- Multi-class SVM: ** $\min_{w} \frac{1}{2} \sum_{k} w_{k}^{T} w_{k} + C \sum_{i} (\max_{k} (\Delta(y_{i}, k) + w_{k}^{T} x - w_{y_{i}}^{T} x))$
- Log-linear model w/ MAP (multi-class) $\min_{k \in \{1,2,...,K\}} \frac{1}{2} \sum_{k} w_{k}^{T} w_{k} + C \sum_{i} [\log_{k \in \{1,2,...,K\}} \exp(w_{k}^{T} x_{i}) - w_{y_{i}}^{T} x_{i}]$

$$\min_{w} \ \frac{1}{2} w^{T} w + C \sum_{i} \log(1 + e^{-\frac{1}{2}})$$

