

# SVMs: Non-Separable Data, Convex Surrogate Loss, Multi-Class Classification, Kernels

Karl Stratos

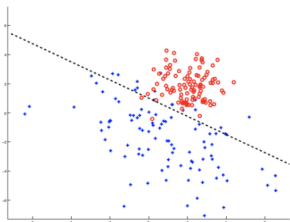
June 21, 2018

## Tangent: Some Loose Ends in Logistic Regression

- ▶ Polynomial feature expansion in logistic regression
- ▶ Regularization in logistic regression
- ▶ Classification metrics

# Increasing the Complexity of Logistic Regression

- ▶ The predicted probability of a logistic regressor is  $p(1|\mathbf{x}, \mathbf{w}, w_0) = \sigma(\mathbf{w} \cdot \mathbf{x} + w_0)$ .
- ▶ Linear decision boundary



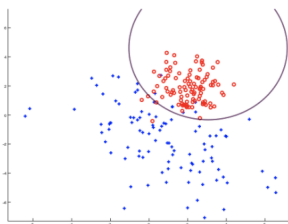
$$\mathbf{w} \cdot \mathbf{x} + w_0 = 0 \quad \Leftrightarrow \quad p(1|\cdot) = \frac{1}{2}$$

## Polynomial Feature Expansion

- ▶ We can use the same polynomial feature expansion that we used in linear regression.
- ▶ For instance, with  $d = 2$  dimensions

$$p(1|\mathbf{x}, \mathbf{w}, w_0) = \sigma(w_0 + w_1x_1 + w_2x_2 + w_3x_1^2 + w_4x_2^2)$$

- ▶ Nonlinear decision boundary



$$w_0 + w_1x_1 + w_2x_2 + w_3x_1^2 + w_4x_2^2 = 0 \quad \Leftrightarrow \quad p(1|\cdot) = \frac{1}{2}$$

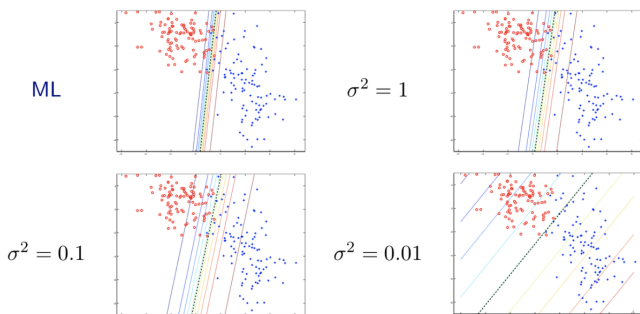
## Tangent: Some Loose Ends in Logistic Regression

- ▶ Polynomial feature expansion in logistic regression
- ▶ Regularization in logistic regression
- ▶ Classification metrics

## Regularization in Logistic Regression

- ▶ Same idea as in linear regression: penalize the squared  $l_2$  or  $l_1$  norm of the model parameter to prevent the model from becoming too “confident” about the training data.
- ▶ Squared  $l_2$  regularization with hyperparameter  $\sigma^2$

$$-\sum_i \log p(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) + \frac{1}{2\sigma^2} \|\mathbf{w}\|^2$$



## Tangent: Some Loose Ends in Logistic Regression

- ▶ Polynomial feature expansion in logistic regression
- ▶ Regularization in logistic regression
- ▶ Classification metrics

## Label Imbalance Problem

- ▶ Suppose most of the labels are  $y = 0$ , say 99.9% of the time.
- ▶ In this scenario, classification accuracy is not a useful metric.

$$\frac{tp + tn}{tp + fp + tn + fn}$$

- ▶ Just by guessing 0 all the time, we get accuracy 99.9%!



## Label Imbalance Problem

- ▶ Suppose most of the labels are  $y = 0$ , say 99.9% of the time.
- ▶ In this scenario, classification accuracy is not a useful metric.

$$\frac{tp + tn}{tp + fp + tn + fn}$$

- ▶ Just by guessing 0 all the time, we get accuracy 99.9%!
- ▶ Consider other metrics, such as
  - ▶ **Precision:** Out of your  $y = 1$  predictions, how many were actually 1?

$$\frac{tp}{tp + fp}$$

- ▶ **Recall:** Out of points that are labeled 1, how many did you label as  $y = 1$ ?

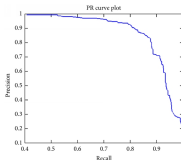
$$\frac{tp}{tp + fn}$$

## More Metrics

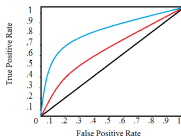
- ▶ **F1**: Harmonic mean of precision and recall

$$2 \frac{\textit{precision} \cdot \textit{recall}}{\textit{precision} + \textit{recall}}$$

- ▶ **PR curve**: change threshold and plot precision/recall



- ▶ **ROC curve**: change threshold and plot false/true positive rates



Back to SVMs

## Review: Basic Support Vector Machines (SVMs)

- ▶ Training data:  $S = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$  where  $y^{(i)} \in \{\pm 1\}$  is a *binary* label of  $\mathbf{x}^{(i)} \in \mathbb{R}^d$ .
- ▶  $S$  is assumed to be *linearly separable*: there exists  $\mathbf{w}$  such that  $y^{(i)}\mathbf{w} \cdot \mathbf{x}^{(i)} > 0$  for all  $i = 1 \dots n$ .
- ▶ Find a separator that maximizes the **margin** on  $S$ :

## Review: Basic Support Vector Machines (SVMs)

- ▶ Training data:  $S = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$  where  $y^{(i)} \in \{\pm 1\}$  is a *binary* label of  $\mathbf{x}^{(i)} \in \mathbb{R}^d$ .
- ▶  $S$  is assumed to be *linearly separable*: there exists  $\mathbf{w}$  such that  $y^{(i)}\mathbf{w} \cdot \mathbf{x}^{(i)} > 0$  for all  $i = 1 \dots n$ .
- ▶ Find a separator that maximizes the **margin** on  $S$ :

$$\begin{aligned}\mathbf{w}_S^{\text{svm}} &:= \arg \max_{\mathbf{w} \in \mathbb{R}^d} \min_{i=1}^n y^{(i)} \frac{\mathbf{w}}{\|\mathbf{w}\|} \cdot \mathbf{x}^{(i)} \\ &\stackrel{\text{wlog}}{\equiv} \arg \max_{\mathbf{w} \in \mathbb{R}^d: \|\mathbf{w}\|=1} \min_{i=1}^n y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)} \\ &\propto \arg \min_{\substack{\mathbf{w} \in \mathbb{R}^d: \\ y^{(i)}\mathbf{w} \cdot \mathbf{x}^{(i)} \geq 1 \forall i}} \|\mathbf{w}\|^2\end{aligned}$$

... Equivalently, find a separator with the minimum  $l_2$  norm.

## Review: Inner Product Formulation

Using the **representer theorem**

$$\exists \beta_1 \dots \beta_n \in \mathbb{R} : \mathbf{w}_S^{\text{svm}} = \sum_{i=1}^n \beta_i \mathbf{x}^{(i)}$$

we can convert the original problem into an equivalent problem

$$\begin{aligned} \min_{\beta_1 \dots \beta_n \in \mathbb{R}} \quad & \sum_{i,j=1}^n \beta_i \beta_j \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)} \\ \text{subject to} \quad & y^{(i)} \sum_{j=1}^n \beta_j \mathbf{x}^{(j)} \cdot \mathbf{x}^{(i)} \geq 1 \quad \forall i = 1 \dots n \end{aligned}$$

where the only information from data we need for training is the **inner product between input points**.

- ▶ Likewise at test time:  $\mathbf{w}_S^{\text{svm}} \cdot \mathbf{x} = \sum_{i=1: \beta_i \neq 0}^n \beta_i \mathbf{x}^{(i)} \cdot \mathbf{x}$
- ▶ Allows for the use of **kernels** (later).

# Today

- ▶ How to handle **non-separable data**
  - ▶ Connection to a convex surrogate loss on the 0-1 loss
- ▶ How to handle **multi-class classification**
- ▶ The **kernel trick**

# Overview

Non-Separable Data

Multi-Class Classification

Kernel Trick



## Introduce Slack Variables

$$\begin{aligned} & \min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{w}\|^2 \\ & \text{subject to } y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)} \geq 1 \quad \forall i = 1 \dots n \end{aligned}$$

⇓

$$\begin{aligned} & \min_{\mathbf{w} \in \mathbb{R}^d, \xi_1 \dots \xi_n \in \mathbb{R}} \|\mathbf{w}\|^2 + \sum_{i=1}^n \xi_i \\ & \text{subject to } y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)} \geq 1 - \xi_i \quad \forall i = 1 \dots n \\ & \quad \quad \quad \xi_i \geq 0 \quad \forall i = 1 \dots n \end{aligned}$$

# Unconstrained Formulation

“Soft” SVM solution

$$\mathbf{w}_S^{\text{soft}}, \xi_1^* \dots \xi_n^* := \arg \min_{\mathbf{w} \in \mathbb{R}^d, \xi_1 \dots \xi_n \in \mathbb{R}} \|\mathbf{w}\|^2 + \sum_{i=1}^n \xi_i$$

with constraints  $\xi_i \geq \max(0, 1 - y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)})$  for all  $i = 1 \dots n$ .

Note that  $\xi_i^* = \max(0, 1 - y^{(i)} \mathbf{w}_S^{\text{soft}} \cdot \mathbf{x}^{(i)})$ , so

$$\mathbf{w}_S^{\text{soft}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{w}\|^2 + \sum_{i=1}^n \max(0, 1 - y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)})$$

No constraints :) Convex but not differentiable, can still be optimized by subgradient descent

## Soft SVMs as Empirical Risk Minimization

- ▶ What we really want: minimize the 0-1 loss

$$\mathbf{w}_S^* = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^n \left[ \left[ y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)} \leq 0 \right] \right]$$

where  $[[\tau]]$  is 1 if  $\tau$  is true and 0 otherwise

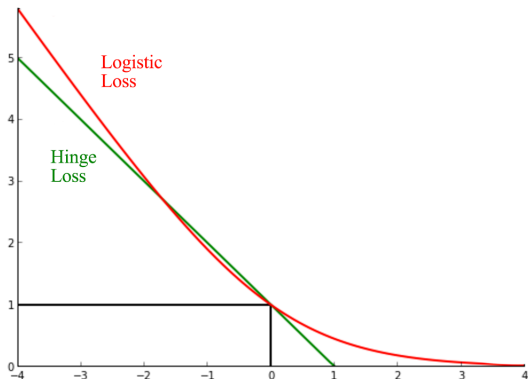
- ▶ Difficult to optimize (neither convex nor differentiable)
- ▶ Instead minimize the **hinge loss**

$$\mathbf{w}_S^* = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^n \underbrace{\max \left( 0, 1 - y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)} \right)}_{\text{hinge}(\mathbf{w} \cdot \mathbf{x}^{(i)})}$$

which is a convex upper bound on the 0-1 loss

# A Big Picture of Binary Classification

Both logistic regression and soft SVMs are  $l_2$ -regularized minimization of the 0-1 loss by convex surrogates.



## Generalized Representer Theorem

**Claim.** Let  $l : \mathbb{R}^n \rightarrow [0, \infty)$  be any function and define

$$\mathbf{w}^* = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{w}\|^2 + l(\mathbf{w} \cdot \mathbf{x}^{(1)}, \dots, \mathbf{w} \cdot \mathbf{x}^{(n)})$$

Then  $\mathbf{w}^* = \sum_{i=1}^n \beta_i \mathbf{x}^{(i)}$  for some  $\beta_1 \dots \beta_n \in \mathbb{R}$ .

**Proof.** Same as in the hard SVM case

Thus we can similarly derive an inner product formulation of the soft SVM solution (will come back to this later):

$$\mathbf{w}_S^{\text{soft}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{w}\|^2 + \underbrace{\sum_{i=1}^n \max(0, 1 - y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)})}_{l(\mathbf{w} \cdot \mathbf{x}^{(1)}, \dots, \mathbf{w} \cdot \mathbf{x}^{(n)})}$$

# Overview

Non-Separable Data

Multi-Class Classification

Kernel Trick

## One-Vs-All

- ▶ Training data:  $S = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$  where  $y^{(i)} \in \{1 \dots m\}$  is the label of  $\mathbf{x}^{(i)} \in \mathbb{R}^d$ .
- ▶ Parameters:  $\mathbf{w}^y \in \mathbb{R}^d$  for each  $y \in \{1 \dots m\}$
- ▶ “One-vs-all” soft SVM objective: optimize

$$\min_{\mathbf{w} \in \mathbb{R}^d, \xi_1 \dots \xi_n \in \mathbb{R}} \|\mathbf{w}\|^2 + \sum_{i=1}^n \xi_i$$

such that  $\xi_i \geq 0$  for all  $i = 1 \dots n$  and

$$\mathbf{w}^{y^{(i)}} \cdot \mathbf{x}^{(i)} - \mathbf{w}^y \cdot \mathbf{x}^{(i)} \geq 1 - \xi_i$$

for all  $i = 1 \dots n$  and  $y \in \{1 \dots m\}$  such that  $y \neq y^{(i)}$

# Unconstrained Formulation

$$\mathbf{w}_S^{\text{one-vs-all}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{w}\|^2 + \sum_{i=1}^n \sum_{y \in \{1 \dots m\}: y \neq y^{(i)}} \max\left(0, 1 + \mathbf{w}^y \cdot \mathbf{x}^{(i)} - \mathbf{w}^{y^{(i)}} \cdot \mathbf{x}^{(i)}\right)$$



# One-Vs-One

- ▶ “One-vs-one” soft SVM objective: optimize

$$\min_{\mathbf{w} \in \mathbb{R}^d, \xi_1 \dots \xi_n \in \mathbb{R}} \|\mathbf{w}\|^2 + \sum_{i=1}^n \xi_i$$

such that  $\xi_i \geq 0$  for all  $i = 1 \dots n$  and

$$\mathbf{w}^{y^{(i)}} \cdot \mathbf{x}^{(i)} - \max_{\mathbf{y} \in \{1 \dots m\}: \mathbf{y} \neq y^{(i)}} \mathbf{w}^{\mathbf{y}} \cdot \mathbf{x}^{(i)} \geq 1 - \xi_i$$

for all  $i = 1 \dots n$

- ▶ Unconstrained Formulation

$$\mathbf{w}_S^{\text{one-vs-all}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{w}\|^2 + \sum_{i=1}^n \max \left( 0, 1 + \max_{\mathbf{y} \in \{1 \dots m\}: \mathbf{y} \neq y^{(i)}} \mathbf{w}^{\mathbf{y}} \cdot \mathbf{x}^{(i)} - \mathbf{w}^{y^{(i)}} \cdot \mathbf{x}^{(i)} \right)$$

# Overview

Non-Separable Data  
Multi-Class Classification  
Kernel Trick

## Inner Product Formulation of Soft SVM (Binary)

- ▶  $G \in \mathbb{R}^{n \times n}$ : a symmetric matrix with  $G_{i,j} = \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)}$  (i.e., the Gram matrix).

$$\boldsymbol{\beta}^* = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^n} \boldsymbol{\beta}^\top G \boldsymbol{\beta} + \sum_{i=1}^n \max \left( 0, 1 - y^{(i)} \boldsymbol{\beta}^\top G_i \right)$$

$$\mathbf{w}_S^{\text{soft}} = \sum_{i=1}^n \beta_i^* \mathbf{x}^{(i)}$$

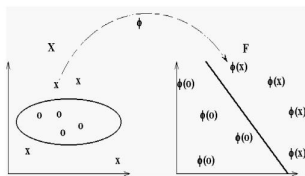
- ▶ User manual

1. Calculate  $G_{i,j} = \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)}$  for every  $i, j = 1 \dots n$ .
2. Find  $\boldsymbol{\beta}^*$  using  $G$  (e.g., by subgradient descent).
3. Test time: given a new point  $\mathbf{x}$  to classify, return

$$\text{sign} \left( \sum_{i=1: \beta_i^* \neq 0}^n \beta_i^* \mathbf{x}^{(i)} \cdot \mathbf{x} \right)$$

## Recall: Polynomial Feature Expansion

- ▶ Idea: transform input by  $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^D$  to allow the linear model to better fit the data.



- ▶ Example: expansion by a degree  $p = 2$  polynomial with bias  $c = 1$

$$\phi^{\text{poly}(2,1)}(x_1 \dots x_d) = \left( (x_i^2)_i, (\sqrt{2}x_i x_j)_{i < j}, (\sqrt{2}x_i)_i, 1 \right)$$

- ▶ **Computationally expensive:** time to calculate feature expansion  $O(d^p)$  exponential in  $p$

## But Computing Inner Product is Easy!

- ▶ Inner product between two points  $\mathbf{x}$  and  $\mathbf{y}$  in the feature space

$$\begin{aligned}\phi^{\text{poly}(2,1)}(\mathbf{x}) \cdot \phi^{\text{poly}(2,1)}(\mathbf{y}) &= \sum_{i,j} x_i x_j y_i y_j + 2 \sum_i x_i y_i + 1 \\ &= \left(\mathbf{x}^\top \mathbf{y} + 1\right)^2\end{aligned}$$

- ▶ Instead of computing  $O(d^2)$  terms in each  $\phi^{\text{poly}(2,1)}(\mathbf{x})$  and  $\phi^{\text{poly}(2,1)}(\mathbf{y})$  and then taking a dot product, we can just
  1. Compute  $z = \mathbf{x}^\top \mathbf{y}$  ( $O(d)$ -time operation)
  2. Square  $z + 1$  ( $O(1)$ -time operation)

## Kernel Trick

- ▶ Idea: when all we need is inner product, we can do “implicit” feature expansion by a kernel function without ever computing the explicit feature expansion
- ▶ **Kernel function**  $K(\mathbf{x}, \mathbf{y})$  is any function that defines pairwise similarity between two data points such that

$$K(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x}) \cdot \phi(\mathbf{y})$$

for *some* input mapping  $\phi : \mathbb{R}^d \rightarrow ?$

- ▶ Applicable beyond SVMs (e.g., kernel PCA)

## Degree- $p$ Polynomial Kernel

Parameters: degree  $p$ , bias  $c$

$$K^{\text{poly}(p,c)}(\mathbf{x}, \mathbf{y}) = \left( \mathbf{x}^\top \mathbf{y} + c \right)^p$$

We saw that the underlying feature expansion is some degree  $p$  polynomial

# Radial Basis Function (RBF) Kernel

Parameter:  $\sigma^2 > 0$

$$K^{\text{RBF}(\sigma^2)}(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{y}\|^2}{2\sigma^2}\right)$$

What is the underlying feature expansion? For  $\sigma^2 = 1$ ,

$$\begin{aligned} K^{\text{RBF}(\sigma^2)}(\mathbf{x}, \mathbf{y}) &= C \sum_{p=0}^{\infty} \frac{1}{p!} (\mathbf{x}^\top \mathbf{y})^p \\ &= C \sum_{p=0}^{\infty} \frac{1}{p!} \phi^{\text{poly}(p,0)}(\mathbf{x}) \cdot \phi^{\text{poly}(p,0)}(\mathbf{y}) \end{aligned}$$

The underlying feature space is *infinite-dimensional*.



# Summary of Kernel Trick for Soft SVM

► Before (“linear kernel”)

1. Calculate  $G_{i,j} = \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)}$  for every  $i, j = 1 \dots n$ .
2. Find  $\beta^*$  using  $G$  (e.g., by subgradient descent).
3. Test time: given a new point  $\mathbf{x}$  to classify, return

$$\text{sign} \left( \sum_{i=1:\beta_i^* \neq 0}^n \beta_i^* \mathbf{x}^{(i)} \cdot \mathbf{x} \right)$$

► Choose some kernel  $K(\mathbf{x}, \mathbf{y})$ .

1. Calculate  $G_{i,j} = K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$  for every  $i, j = 1 \dots n$ .
2. Find  $\beta^*$  using  $G$  (e.g., by subgradient descent).
3. Test time: given a new point  $\mathbf{x}$  to classify, return

$$\text{sign} \left( \sum_{i=1:\beta_i^* \neq 0}^n \beta_i^* K(\mathbf{x}^{(i)}, \mathbf{x}) \right)$$

## Aside: Kernel Approximation

- ▶ Kernel trick is clever but requires the inner product formulation. This requires storing the  $n \times n$  Gram matrix: not scalable.
- ▶ Kernel approximation: approximate the implicit feature expansion defined under a kernel, and use that expansion directly
- ▶ Rahimi and Recht (2007):  $z(\mathbf{x}) \in \mathbb{R}^N$  where  $z_i(\mathbf{x}) := \sqrt{2/N} \cos(\mu_i \cdot \mathbf{x} + b_i)$  given by  $\mu_i \sim \mathcal{N}(0, I_d)$  and  $b_i \sim \mathcal{U}(0, 2\pi)$

$$\mathbf{E}[z(\mathbf{x}) \cdot z(\mathbf{y})] = K^{\text{RBF}(1)}(\mathbf{x}, \mathbf{y})$$

Use  $z(\mathbf{x}) \in \mathbb{R}^N$  directly (no kernel trick)

# Summary

- ▶ Soft SVM = hard SVM + slack variables to handle non-separable data
  - ▶ A unifying framework for SVMs and logistic regression: convex surrogate loss on the 0-1 loss
- ▶ SVMs can be naturally extended to handle multi-class classification
- ▶ Kernel trick: when training and inference only depend on inner product between data points, we can replace the inner product with a kernel function and perform implicit feature expansion