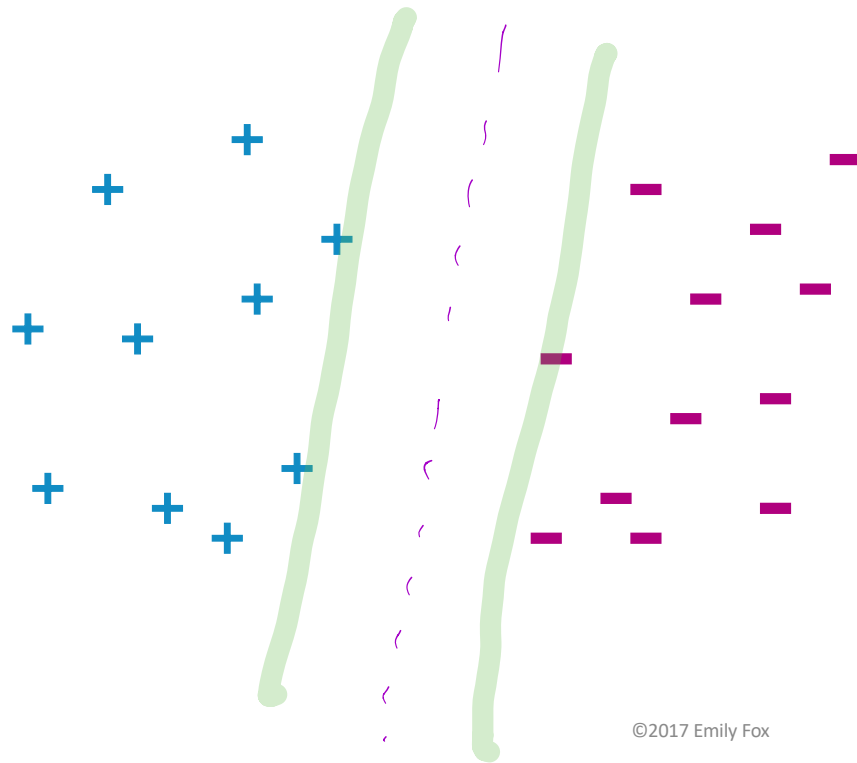


Support Vector Machines

CSE 446: Machine Learning
Slides by Emily Fox + Kevin Jamieson + others
Presented by Anna Karlin

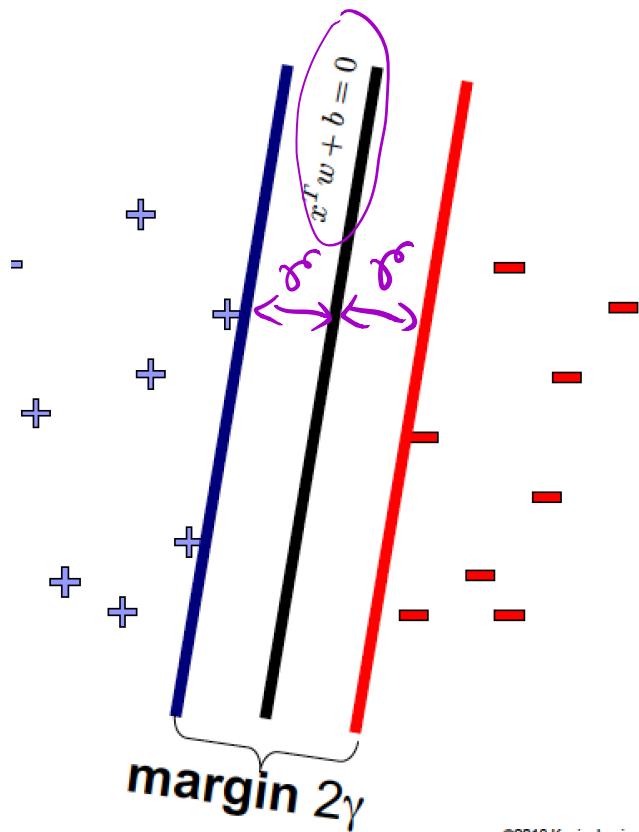
May 1, 2019

Linear classifiers—Which line is better?



$$\{(x_i, y_i)\} \quad x_i \in \mathbb{R}^d \\ y_i \in \{-1, +1\}$$

find linear decision
boundary
with the largest
margin possible



find w & b
that gives largest γ

Maximizing the margin for linearly separable data

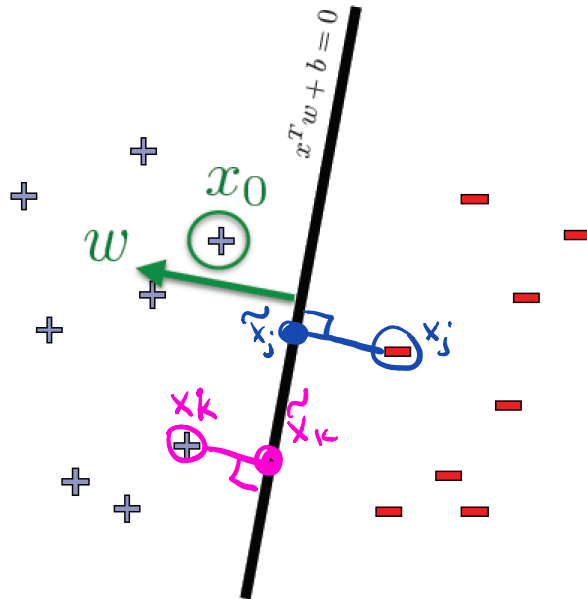
$$\frac{w}{\|w\|} \quad \text{unit vector}$$

$$u \cdot v = \|u\| \|v\| \cos \theta$$

$$\tilde{x}_k \cdot w + b = 0$$

$$\tilde{x}_j \cdot w + b = 0$$

Given hyperplane, what is margin?



Distance from x_0 to
hyperplane defined
by $x^T w + b = 0$?

$$y_k = +1 \quad \|x_k - \tilde{x}_k\| = (x_k - \tilde{x}_k) \cdot \frac{w}{\|w\|} = \frac{x_k \cdot w + b}{\|w\|} \geq \gamma$$

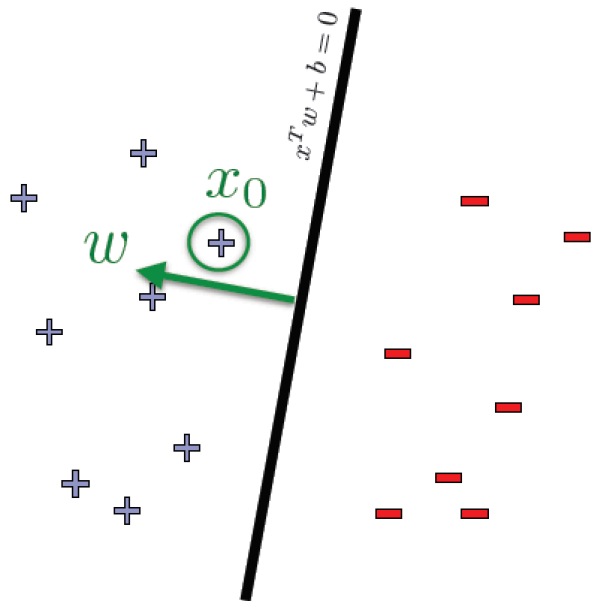
$$y_j = -1 \quad (x_j - \tilde{x}_j) \cdot \frac{w}{\|w\|} = -\|x_j - \tilde{x}_j\|$$

$$\text{want} \quad -\frac{(x_j \cdot w + b)}{\|w\|} \geq \gamma$$

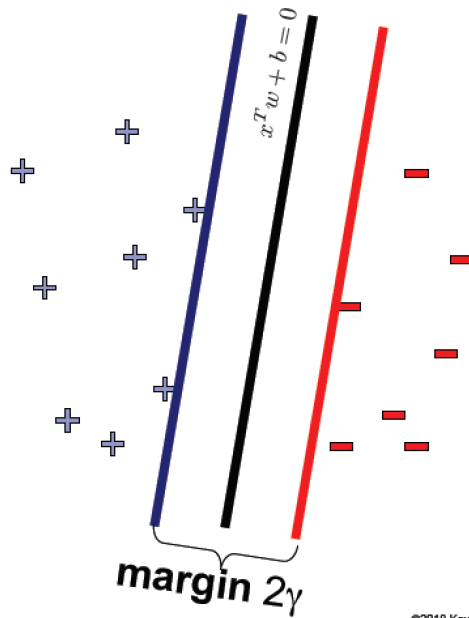
$$\forall 1 \leq i \leq n$$

$$y_i \left(\frac{x_i \cdot w + b}{\|w\|} \right) \geq \gamma$$

Given hyperplane, what is margin?



Our optimization problem



Optimal Hyperplane

$$\max_{w,b} \gamma$$

$$\text{subject to } \underbrace{\frac{1}{\|w\|_2}}_{1/\|w\|_2} \underbrace{y_i(x_i^T w + b)}_{1/\|w\|_2} \geq \gamma \quad \forall i$$

$$y_i(x_i^T w + b) \geq \gamma^0 \|w\|_2$$

\downarrow \uparrow \uparrow
 $1/\|w\|_2$ $1/\|w\|_2$ $1/\|w\|_2$

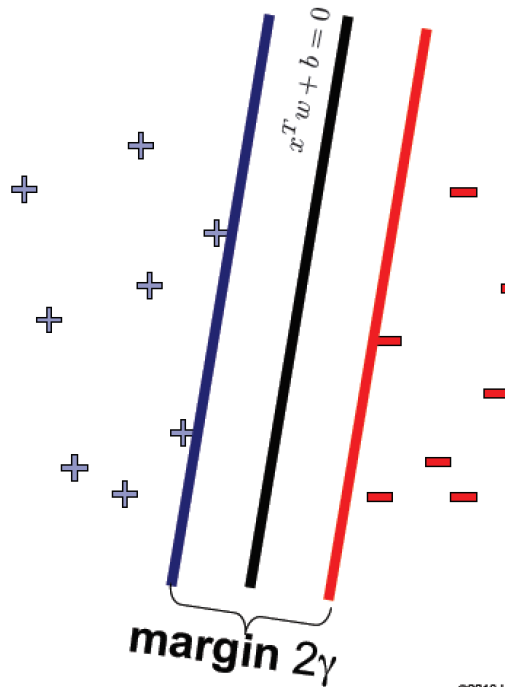
Choose it so that

$$\gamma^0 \|w\|_2 = 1$$

$$\begin{aligned}
 & \max \quad \gamma^0 \\
 & \text{s.t.} \quad y_i(x_i^T w + b) \geq 1 \\
 & \text{where} \quad \gamma^0 = \frac{1}{\|w\|_2}
 \end{aligned}
 \quad \equiv \max \quad \left(\frac{1}{\|w\|_2} \right)
 \quad \equiv \min \|w\|_2
 \quad \equiv \min (\|w\|_2^2)$$

Final version

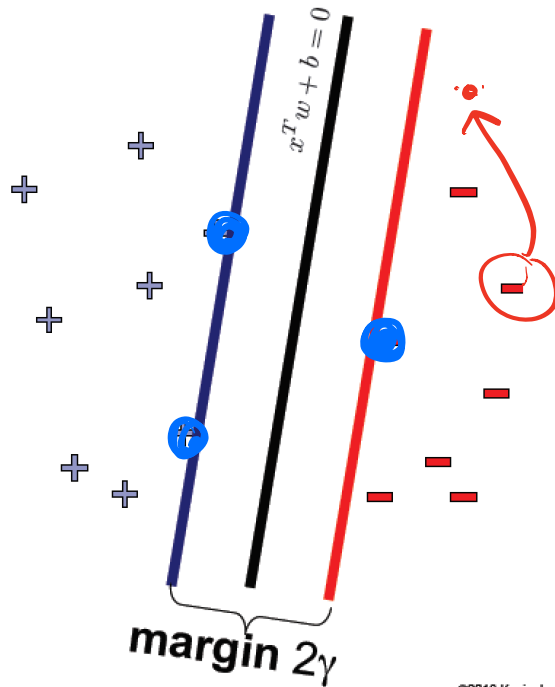
Solvable efficiently –
quadratic programming problem



Optimal Hyperplane (reparameterized)

$$\begin{aligned}
 & \min_{w,b} \|w\|_2^2 \\
 & \text{subject to} \quad y_i(x_i^T w + b) \geq 1 \quad \forall i
 \end{aligned}$$

What are support vectors?



Optimal Hyperplane (reparameterized)

$$\min_{w,b} ||w||_2^2$$

$$\text{subject to } y_i(x_i^T w + b) \geq 1 \quad \forall i$$

What if the data are not linearly separable?

What if data are not linearly separable?

Use feature maps...

$$\phi: \mathbb{R}^d \rightarrow F \quad (\mathbb{R}^p)$$

Solve problem on $\{(\phi(x_i), y_i)\}_{i=1}^n$

SVMs can be kernelized!!!!

write alg in terms of inner products

$$x_i \cdot x_j \Rightarrow \phi(x_i) \cdot \phi(x_j)$$

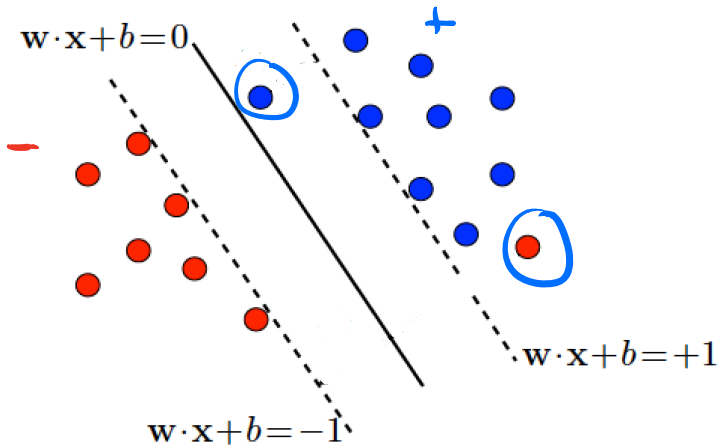
kernel has
an inner product

$$K(x, x') \text{ that implements} \\ = \phi(x) \cdot \phi(x')$$

$$w = \sum_{i=1}^n \alpha_i y_i x_i$$

$\alpha_i \geq 0$

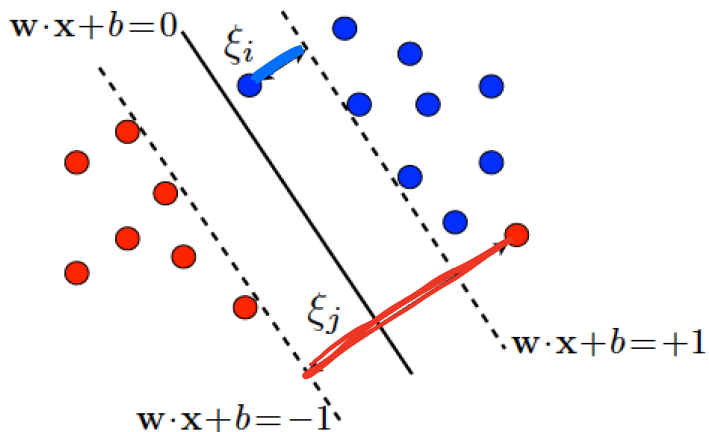
What if data are still not linearly separable?



Courtesy Mehryar Mohri

What if data are still not linearly separable?

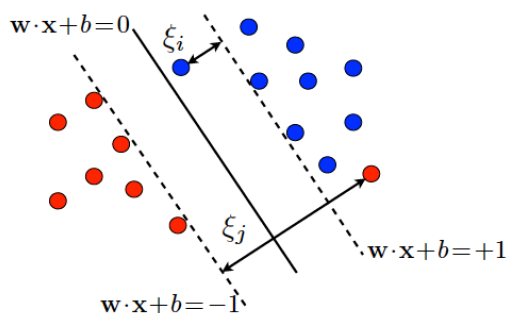
allow for slack



Courtesy Mehryar Mohri

$$\begin{aligned}
 &\min \|w\|_2^2 \\
 &\text{st.} \quad \underline{y_i(w \cdot x_i + b) \geq 1 - \xi_i} \\
 &\quad \xi_i \geq 0 \quad \forall i
 \end{aligned}$$

What if data are still not linearly separable?



$$\min_{\mathbf{w}, b, \xi} \underbrace{\|\mathbf{w}\|^2 + C \sum_i \xi_i}_{\text{total slack}}$$

$$\boxed{y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i} \quad \forall i$$

$$\underline{\xi_i \geq 0} \quad \forall i.$$

Courtesy Mehryar Mohri

Fix (\vec{w}, b)

$$y_j(\mathbf{w} \cdot \mathbf{x}_j + b) \geq 1 \Rightarrow \xi_j = 0$$

$$y_k(\mathbf{w} \cdot \mathbf{x}_k + b) < 1 \Rightarrow \xi_k = 1 - y_k(\mathbf{w} \cdot \mathbf{x}_k + b)$$

CSE 446: Machine Learning

Replace obj fn with

$$\min_{\mathbf{w}, b} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max(0, 1 - y_i(\mathbf{w} \cdot \mathbf{x}_i + b))$$

$$z_+ = \max(z, 0)$$

Final objective



$$\lambda \neq \frac{1}{nC}$$

$$\frac{1}{n} \sum_{i=1}^n (1 - y_i((\mathbf{w}^T \mathbf{x}_i + b)))_+ + \lambda \|\mathbf{w}\|_2^2$$

loss on i^{th}
data pt

hinge loss

Gradient descent for SVMs

0-1 loss on (x, y)

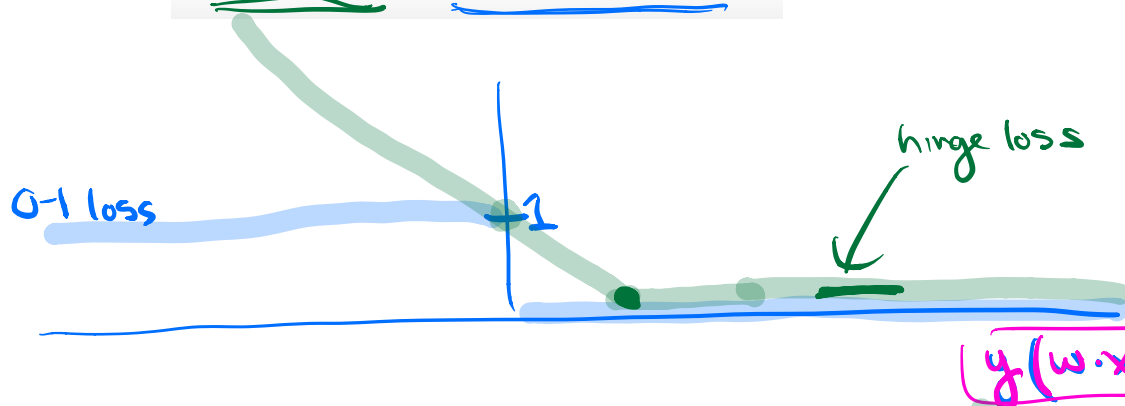
$$1 \cdot \{ \text{sign}(w^T x + b) \neq y \}$$

mistake on this pt

$$\frac{1}{n} \sum_{i=1}^n (1 - y_i ((w^T x_i + b))_+ + \lambda ||w||_2^2$$

Hinge loss

- Hinge loss: $\ell((x, y), w) = (1 - y(w^T x + b))_+$



$y(w \cdot x + b)$
measures correctness
of our prediction

- Subgradient of hinge loss:

$$d=1$$

$$\partial \ell_w = \begin{cases} 0 \\ [-yx, 0] \\ -yx \end{cases}$$

$$y(w \cdot x + b) > 1$$

$$y(w \cdot x + b) = 1$$

$$1 - y(w \cdot x + b) > 0$$

Minimizing regularized hinge loss (aka SVMs)

- Given a dataset: $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)$
- Minimize regularized hinge loss: $\frac{1}{n} \sum_i (1 - y_i((\mathbf{w}^T \mathbf{x}_i + b))_+ + \lambda \|\mathbf{w}\|_2^2$

Subgradient descent for hinge minimization


- Given data: $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)$

- Want to minimize:

$$\frac{1}{n} \sum_{i=1}^n \ell((\mathbf{x}_i, y_i), \mathbf{w}) + \lambda \|\mathbf{w}\|_2^2 = \frac{1}{n} \sum_i (1 - y_i (\mathbf{w}^T \mathbf{x}_i + b))_+ + \lambda \|\mathbf{w}\|_2^2$$

- As we've discussed, subgradient descent works like gradient descent:
 - But if there are multiple subgradients at a point, just pick (any) one:

$$\partial_{\mathbf{w}} \ell((\mathbf{x}, y), \mathbf{w}) = \mathbb{I}\{y(\mathbf{w}^T \mathbf{x} + b) \leq 1\} (-y\mathbf{x})$$

subgradient. \vec{g} is a subgradient at w for $f(\cdot)$
 $f(w') \geq f(w) + \vec{g} \cdot (w' - w)$ 

Subgradient descent for hinge minimization

- Given data: $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)$

- Want to minimize:

$$\frac{1}{n} \sum_{i=1}^n \ell((\mathbf{x}_i, y_i), \mathbf{w}) + \lambda \|\mathbf{w}\|_2^2 = \frac{1}{n} \sum_i (1 - y_i((\mathbf{w}^T \mathbf{x}_i + b))_+ + \lambda \|\mathbf{w}\|_2^2$$

- As we've discussed, subgradient descent works like gradient descent:
 - But if there are multiple subgradients at a point, just pick (any) one:

$$\begin{aligned} \mathbf{w}_{t+1} &:= \mathbf{w}_t - \eta \left(\frac{1}{n} \sum_{i=1}^n \partial_{\mathbf{w}} \ell((\mathbf{x}_i, y_i), \mathbf{w}) + 2\lambda \mathbf{w}_t \right) & \partial_{\mathbf{w}} \ell((\mathbf{x}, y), \mathbf{w}) &= \mathbb{I}\{y(\mathbf{w}^T \mathbf{x} + b) \leq 1\}(-y\mathbf{x}) \\ &= \mathbf{w}_t - \eta \left(\frac{1}{n} \sum_{i=1}^n \mathbb{I}\{y(\mathbf{w}_t \cdot \mathbf{x}_i + b) \leq 1\}(-y_i \mathbf{x}_i) + 2\lambda \mathbf{w}_t \right) \\ &= \mathbf{w}_t + \eta \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{y(\mathbf{w}_t \cdot \mathbf{x}_i + b) \leq 1\}(y_i \mathbf{x}_i) - \eta 2\lambda \mathbf{w}_t. \end{aligned}$$

SVM

- (Sub)gradient Descent Update

$$\begin{aligned}\mathbf{w}_{t+1} &:= \mathbf{w}_t - \eta \left(\frac{1}{n} \sum_{i=1}^n \partial_{\mathbf{w}} \ell((\mathbf{x}_i, y_i), \mathbf{w}) + 2\lambda \mathbf{w}_t \right) \\ &= \mathbf{w}_t - \eta \left(\frac{1}{n} \sum_{i=1}^n \mathbb{I}\{y(\mathbf{w}_t \cdot \mathbf{x}_i + b) \leq 1\} (-y_i \mathbf{x}_i) + 2\lambda \mathbf{w}_t \right) \\ &= \mathbf{w}_t + \eta \left(\frac{1}{n} \sum_{i=1}^n \mathbb{I}\{y(\mathbf{w}_t \cdot \mathbf{x}_i + b) \leq 1\} (y_i \mathbf{x}_i) - \eta 2\lambda \mathbf{w}_t \right).\end{aligned}$$

- SGD update

$$\mathbf{w}_{t+1} := \mathbf{w}_t + \eta \mathbb{I}\{y(\mathbf{w}_t \cdot \mathbf{x}_i + b) \leq 1\} (y_i \mathbf{x}_i) - \eta 2\lambda \mathbf{w}_t.$$

Machine learning problems

- Given i.i.d. data set:

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)$$

- Find parameters \mathbf{w} to minimize average loss
(or regularized version):

$$\frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{w})$$

Squared loss:

$$\ell_i(\mathbf{w}) = (y_i - \mathbf{x}_i^T \mathbf{w})^2$$

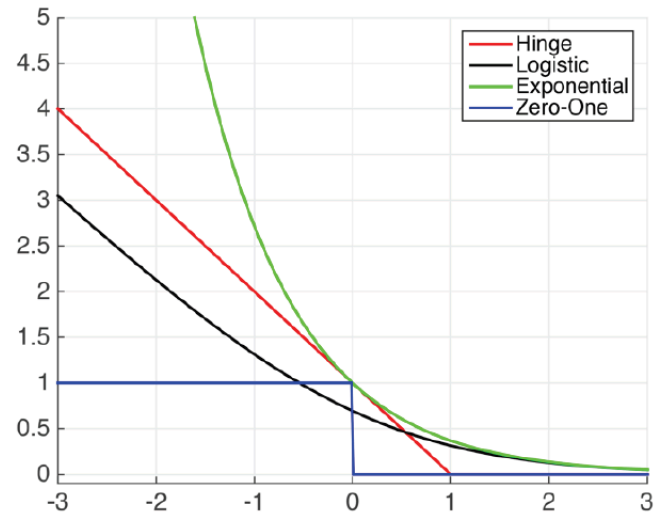
Logistic loss:

$$\ell_i(\mathbf{w}) = \log(1 + \exp(-y_i \mathbf{x}_i^T \mathbf{w}))$$

Hinge loss:

$$\ell_i(\mathbf{w}) = \max\{0, 1 - y_i \mathbf{x}_i^T \mathbf{w}\}$$

Courtesy Killian Weinberger



What you need to know...

- Maximizing margin
- Derivation of SVM formulation
- Non-linearly separable case
 - Hinge loss
 - a.k.a. adding slack variables
- Can optimize SVMs with SGD
 - Many other approaches possible