# SVMs

# Two different approaches to regression/classification

*Generative*

- **Assume something about P(x,y)** $P\left(y \mid x\right)$
- **Find f which maximizes likelihood of training data assumption**
  - **Often reformulated as minimizing loss**

**Versus**

*Discriminative*

- **Pick a loss function**
- **Pick a set of hypotheses H** *linear & NN --*
- **Pick f from H which minimizes loss on training data**

# Our description of logistic regression was the former

- **Learn**: f:$\mathbf{X} \longrightarrow Y$
  - **X** – features
  - **Y – target classes**

$$Y \in \{-1, 1\}$$

- **Expected loss of f:**

$$\mathbb{E}_{XY}[\mathbf{1}\{f(X) \neq Y\}] = \mathbb{E}_X[\mathbb{E}_{Y|X}[\mathbf{1}\{f(x) \neq Y\}|X = x]]$$

$$\mathbb{E}_{Y|X}[\mathbf{1}\{f(x) \neq Y\}|X = x] = 1 - P(Y = f(x)|X = x)$$

- **Bayes optimal classifier:**

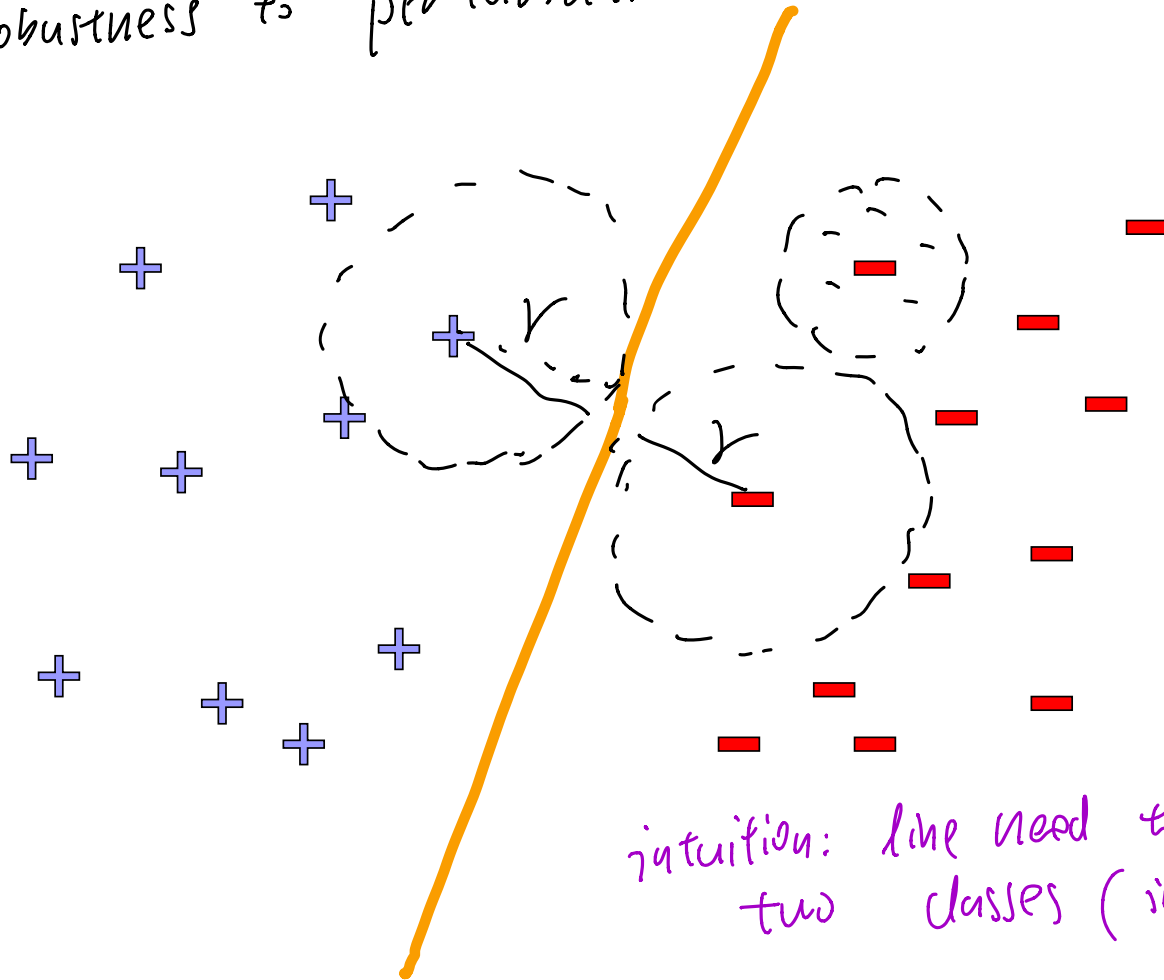$$f(x) = \arg\max_{y} \mathbb{P}(Y = y|X = x)$$

- **Model of logistic regression:**

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

- **Loss function:**

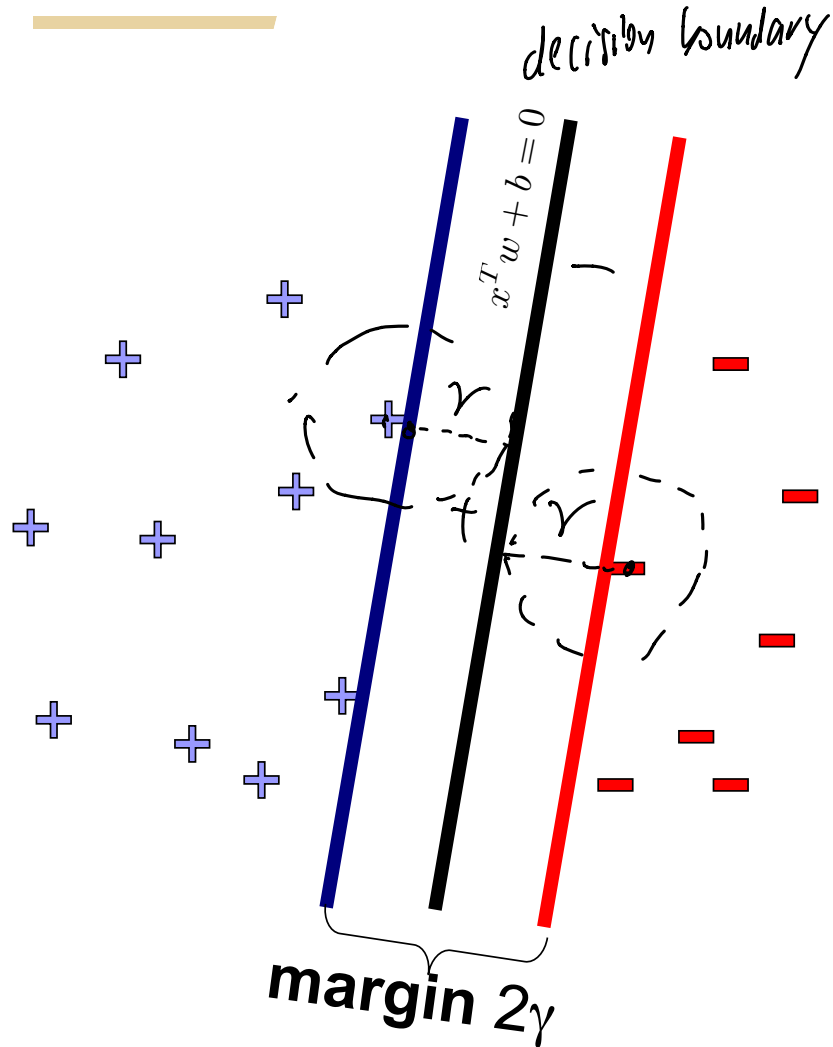$$\ell(f(x), y) = \mathbf{1}\{f(x) \neq y\}$$

**What if the model is wrong? What other ways can we pick linear decision rules?**

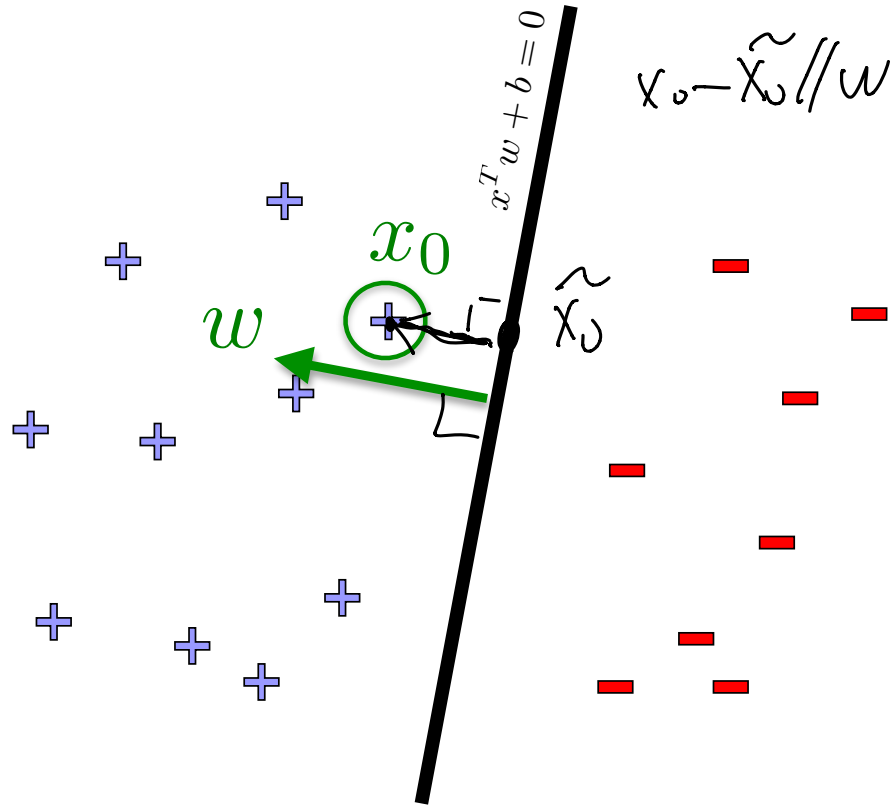# Linear classifiers – Which line is better?

robustness to perturbation

intuition: line need to balance
two classes ( in the middle)

# Pick the one with the largest margin!

decision boundary

$x^T w + b = 0$

$\gamma$

$\gamma$

margin $2\gamma$

$\gamma$: margin

the minimum distance from training point to decision boundary

Goal: $(w, b)$
find decision boundary with maximum $\gamma$

# Pick the one with the largest margin!



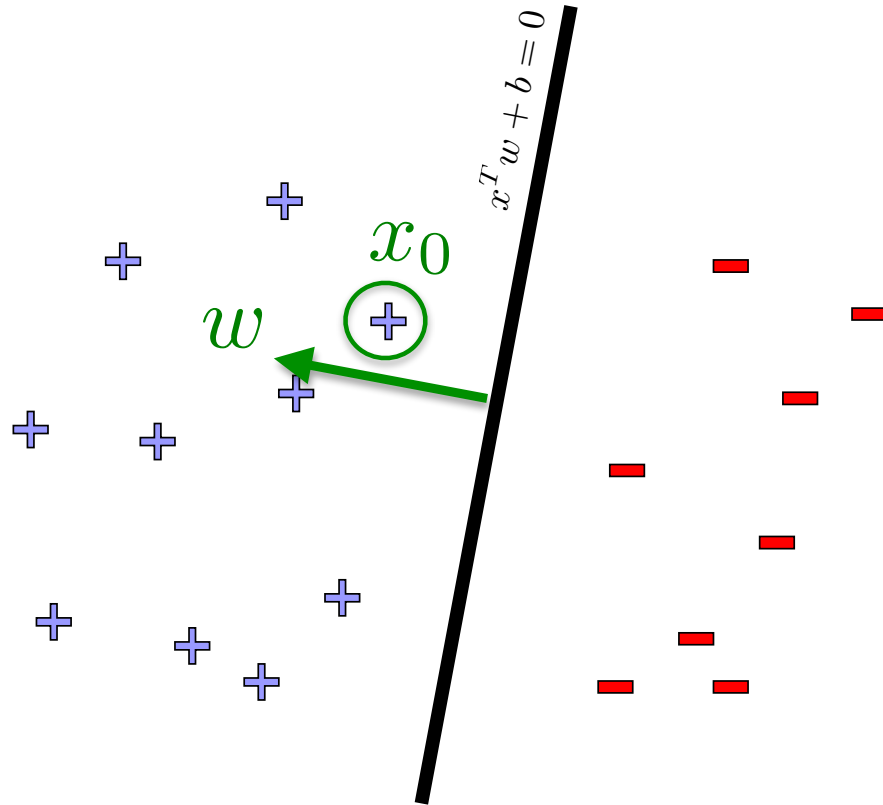$x^T w + b = 0$

$x_0$

$w$

$\tilde{x_0}$

$x_0 - \tilde{x_0} // w$

$\tilde{x_0}^T w + b = 0$ $\left( \begin{array}{l} \tilde{x_0} \text{ is on} \\ \text{decision} \\ \text{boundary} \end{array} \right)$

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

distance :

$\| x_0 - \tilde{x_0} \|_2$

$= \left| (x_0 - \tilde{x})^T \dfrac{w}{\|w\|_2} \right|$ $\left( \begin{array}{l} \text{projection} \\ \text{on its direction} \end{array} \right)$

$= \dfrac{1}{\|w\|_2} \left| x_0^T w - \tilde{x}^T w \right|$

$= \dfrac{1}{\|w\|_2} \left| x_0^T w + b \right|$

# Pick the one with the largest margin!



$x^T w + b = 0$

$x_0$

$w$

Distance from $x_0$ to hyperplane defined by $x^T w + b = 0$?
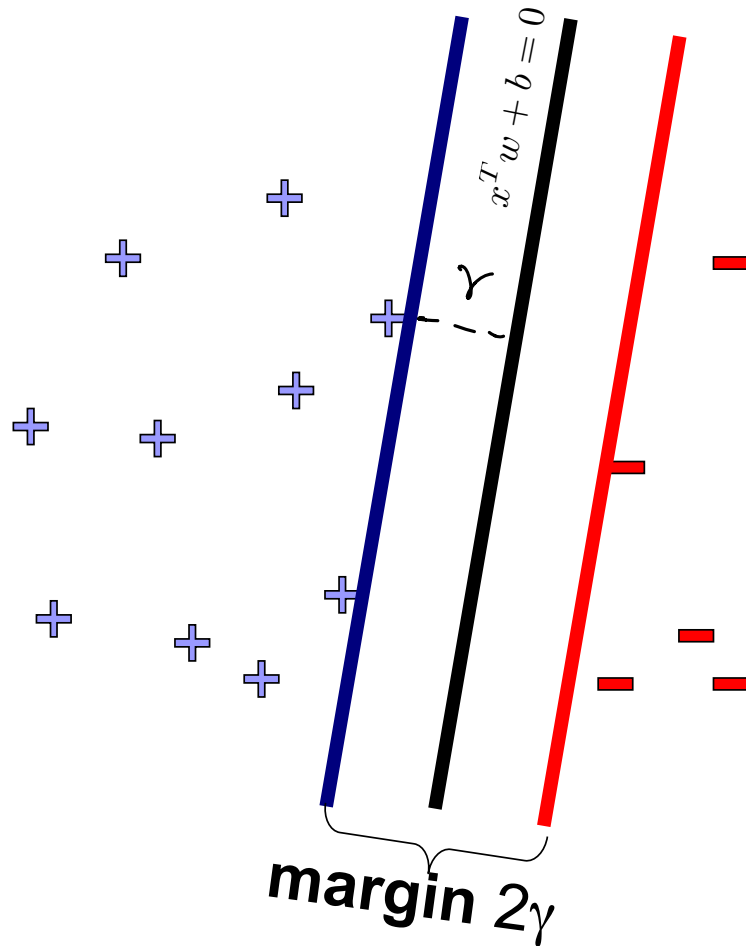
If $\widetilde{x}_0$ is the projection of $x_0$ onto the hyperplane then
$$||x_0 - \widetilde{x}_0||_2 = |(x_0^T - \widetilde{x}_0)^T \frac{w}{||w||_2}|$$

$$= \frac{1}{||w||_2}|x_0^T w - \widetilde{x}_0^T w|$$

$$= \frac{1}{||w||_2}|x_0^T w + b|$$

$\underbrace{\phantom{xxxxxxxx}}_{\text{normalization}}$ $\underbrace{\phantom{xxxxxx}}_{\begin{array}{c}\text{decision}\\\text{rule}\end{array}}$

# Pick the one with the largest margin!



$x^T w + b = 0$

margin $2\gamma$

Distance of $x_0$ from hyperplane $x^T w + b$:

want $\dfrac{1}{||w||_2}(x_0^T w + b) \geq \gamma$

Optimal Hyperplane

$\max\limits_{w, b} \ \gamma$

subject to: $\dfrac{1}{||w||_2} y_i (x_i^T w + b) \geq \gamma, \ \forall i$

$y_i \in \{1, -1\}$

# Pick the one with the largest margin!



$x^T w + b = 0$

margin $2\gamma$

Distance of $x_0$ from hyperplane $x^T w + b$:
$$\frac{1}{||w||_2}(x_0^T w + b)$$

Optimal Hyperplane

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

# Pick the one with the largest margin!

Change of variable:

let $\tilde{w} = \frac{w}{||w||_2}\gamma$   (1)

$\tilde{b} = \frac{b}{||w||_2 \gamma}$

$(1) \Rightarrow ||\tilde{w}||_2 = |\frac{w}{||w||_2}| \cdot \frac{1}{\gamma} \Rightarrow \gamma = \frac{1}{||\tilde{w}||_2}$

Distance of $x_0$ from hyperplane $x^T w + b$:
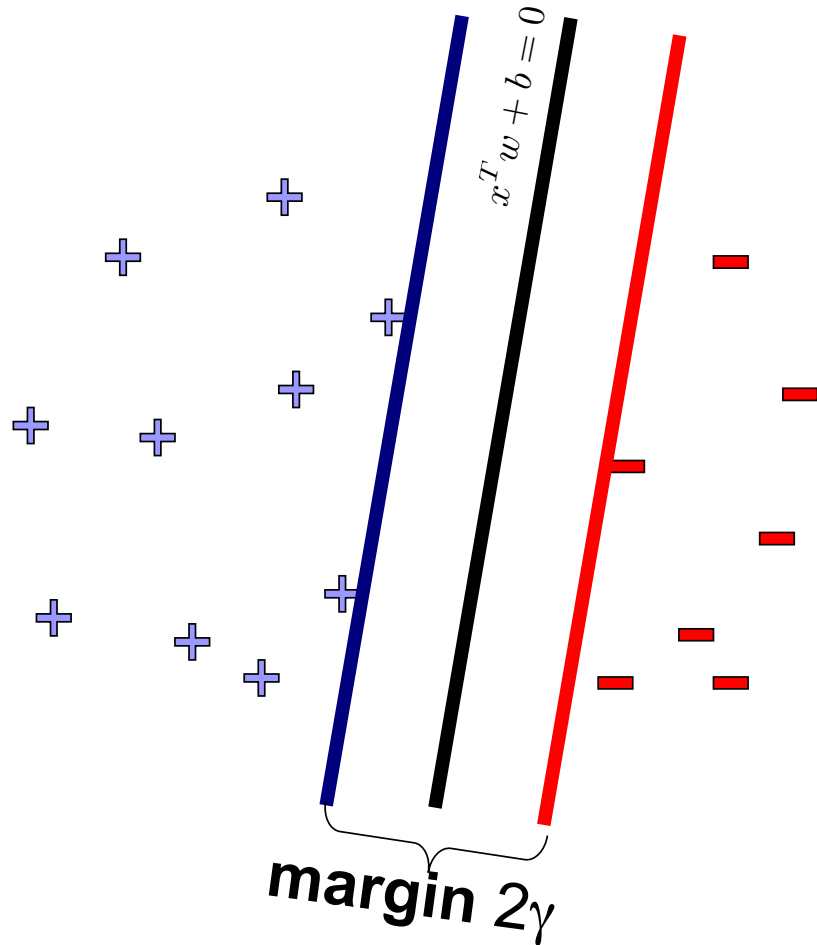
$$\frac{1}{||w||_2}(x_0^T w + b)$$

$x^T w + b = 0$

**margin $2\gamma$**

Optimal Hyperplane

$$\max_{w,b} \gamma \qquad \text{non-convex}$$

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

Optimal Hyperplane (reparameterized)

$\max_{\tilde{w}, \tilde{b}} \frac{1}{||\tilde{w}||_2} \iff \min_{\tilde{w}, \tilde{b}} ||\tilde{w}||_2^2$

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

# Pick the one with the largest margin!



$x^T w + b = 0$

margin $2\gamma$

Distance of $x_0$ from hyperplane $x^T w + b$:

$$\frac{1}{||w||_2}(x_0^T w + b)$$
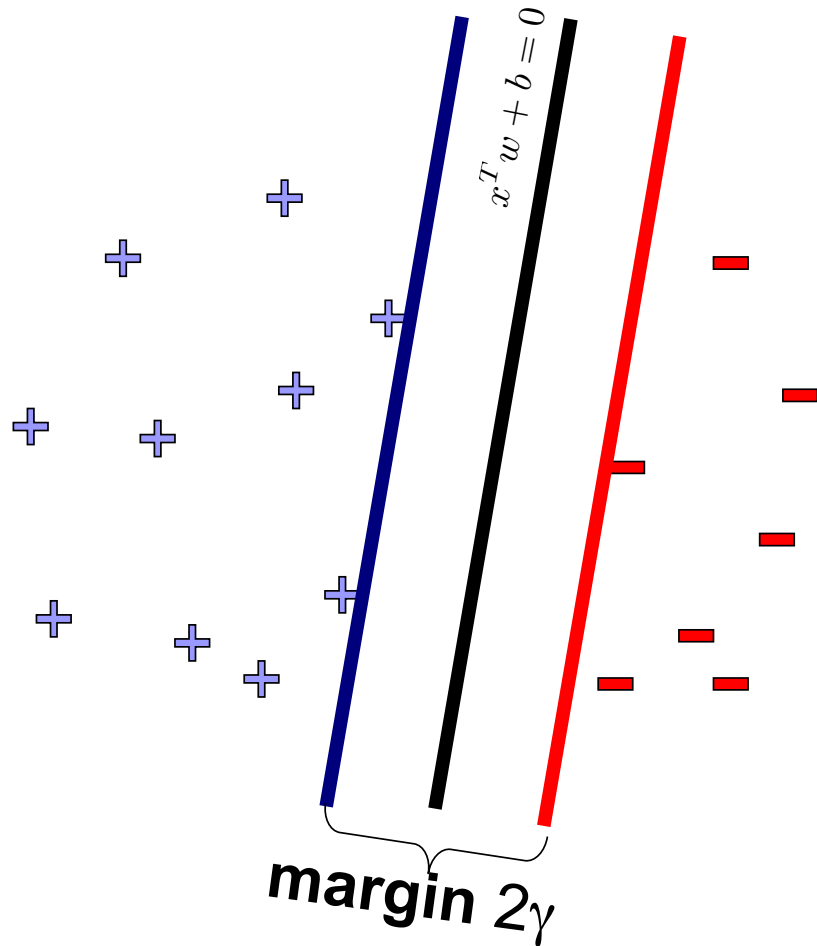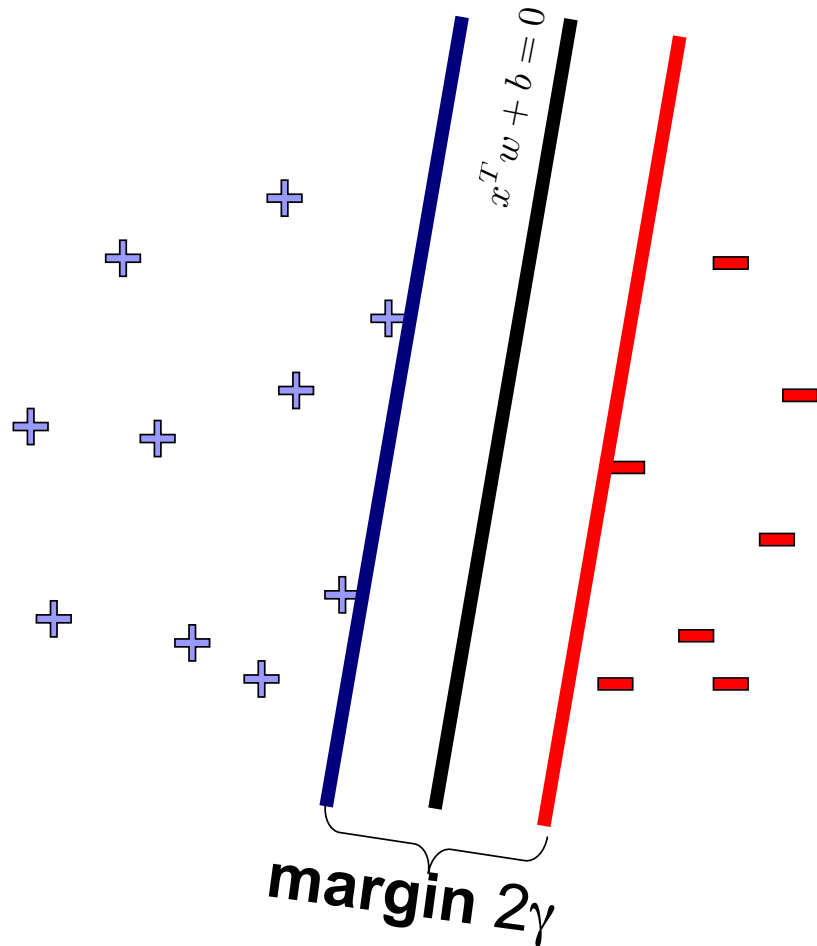
Optimal Hyperplane

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

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

Optimal Hyperplane (reparameterized)

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

$$\text{subject to } y_i(x_i^T w + b) \geq 1 \quad \forall i$$
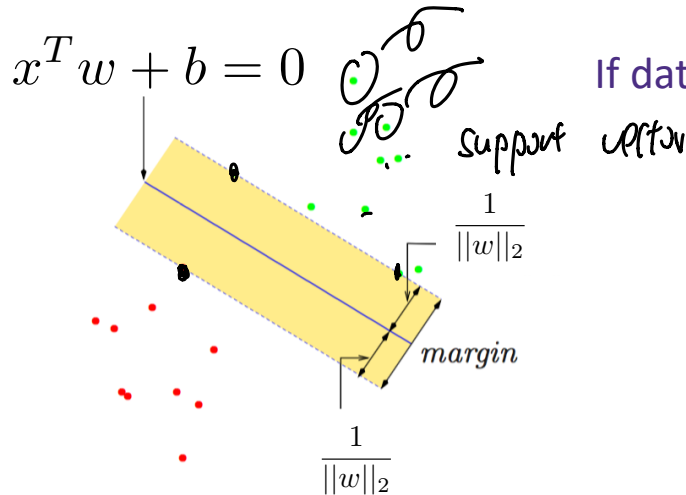
# Pick the one with the largest margin!

$x^T w + b = 0$

margin $2\gamma$

- ■ Solve efficiently by many methods, e.g.,
  - □ quadratic programming (QP)
    - ▪ Well-studied solution algorithms
  - □ Stochastic gradient descent
  - □ Coordinate descent (in the dual)

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 are support vectors

$$x^T w + b = 0$$

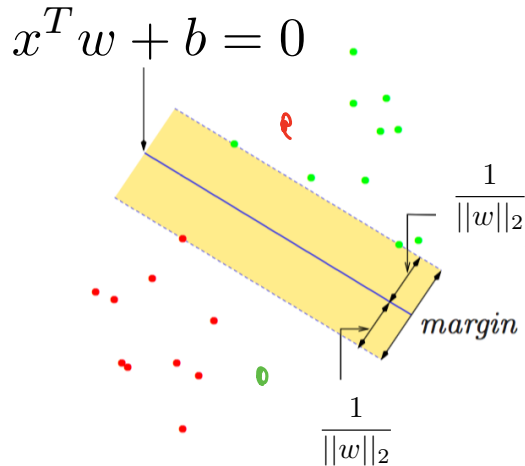*Support vector*

If data is linearly separable

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

$$y_i(x_i^T w + b) \geq 1 \quad \forall i$$

$$x_i, \quad y_i(x_i^T w + b) = 1$$

$$\frac{1}{\|w\|_2}$$

*margin*

$$\frac{1}{\|w\|_2}$$

Note: the solution of this can be written in terms of very few of the training points. These points are known as support vectors.

# What if the data is not linearly separable?

$$x^T w + b = 0$$



$$\frac{1}{||w||_2}$$

*margin*

$$\frac{1}{||w||_2}$$

If data is linearly separable

$$\min_{w,b} ||w||_2^2$$
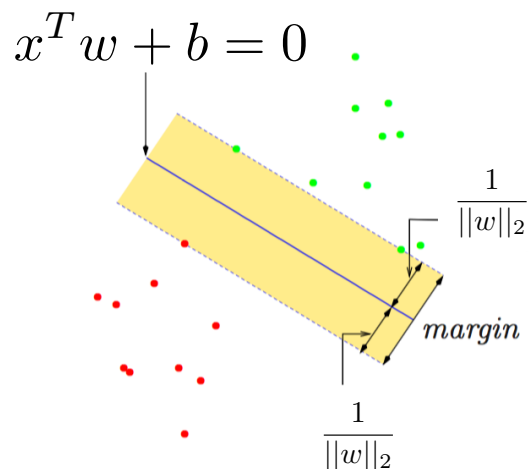
$$y_i(x_i^T w + b) \geq 1 \quad \forall i$$

If data is not linearly separable, some points don't satisfy margin constraint:

Two options:
1. Introduce slack to this optimization problem
2. Lift to higher dimensional space ← kernel

# What if the data is not linearly separable?

$$x^T w + b = 0$$



$$\frac{1}{||w||_2}$$

margin

$$\frac{1}{||w||_2}$$

$$x^T w + b = 0$$

$y_i = -1$

$x^T w + b = -1$

$\xi_4^*$   $\xi_5^*$

$\xi_1^*$   $\xi_3^*$

$\xi_2^*$

$$\frac{1}{||w||_2}$$

$y_i = 1$

$x^T w + b = 1$

margin

$$\frac{1}{||w||_2}$$

If data is linearly separable:

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

$$y_i(x_i^T w + b) \geq 1 \quad \forall i$$

If data is not linearly separable,
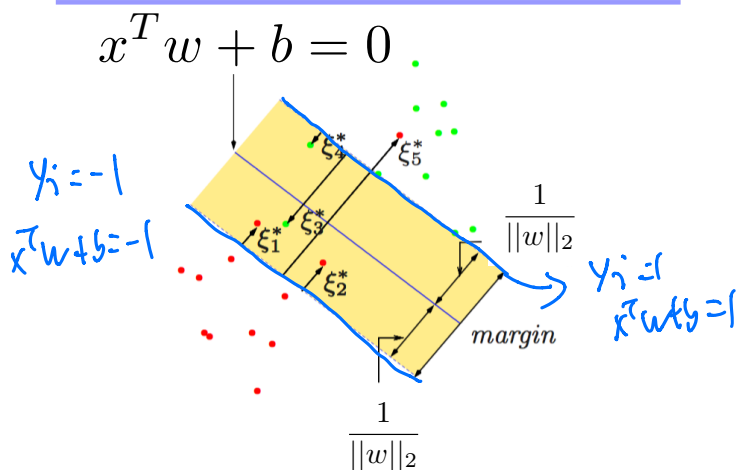some points don't satisfy margin constraint:

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

slack variable

$$y_i(x_i^T w + b) \geq 1 - \xi_i \quad \forall i$$

$$\xi_i \geq 0, \sum_{j=1}^{n} \xi_j \leq \nu$$

convex

# SVM as penalization method

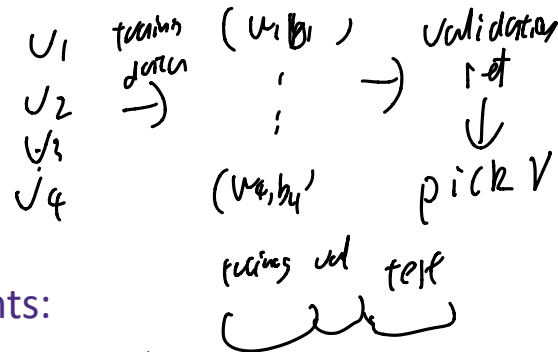- Original quadratic program with linear constraints:

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

$$y_i(x_i^T w + b) \geq 1 - \underline{\xi_i} \quad \forall i$$

$$\xi_i \geq 0, \sum_{j=1}^{n} \xi_j \leq \nu$$

make margin big

violates constraints
as little as possible

# SVM as penalization method

- Original quadratic program with linear constraints:

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

$$y_i(x_i^T w + b) \geq 1 - \xi_i \quad \forall i$$

$$\xi_i \geq 0, \sum_{j=1}^{n} \xi_j \leq \nu$$

KKT condition

- Using same constrained convex optimization trick as for lasso:
  For any $\nu \geq 0$ there exists a $\lambda \geq 0$ such that the solution
  the following solution is equivalent:

$$\min \quad \sum_{i=1}^{n} \max\{0, 1 - y_i(b + x_i^T w)\} + \lambda ||w||_2^2$$

constraints          margin

# SVMs: optimizing what?

SVM objective:

data set: $\{(x_i, y_i)\}_{i=1}^{u}$

$x_i \in \mathbb{R}^d$

$y_i \in \{1, -1\}$

$$\frac{1}{h} \sum_{i=1}^{n} \max\{0, 1 - y_i(b + x_i^T w)\} + \lambda ||w||_2^2 \quad = \sum_{i=1}^{n} \ell_i(w, b)$$

$$\nabla_w \ell_i(w, b) = \begin{cases} -x_i y_i + \frac{2\lambda}{n} w & \text{if } y_i(b + x_i^T w) < 1 \\ \frac{2\lambda}{n} & \text{otherwise} \end{cases}$$

$$\nabla_b \ell_i(w, b) = \begin{cases} -y_i & \text{if } y_i(b + x_i^T w) < 1 \\ 0 & \text{otherwise} \end{cases}$$