



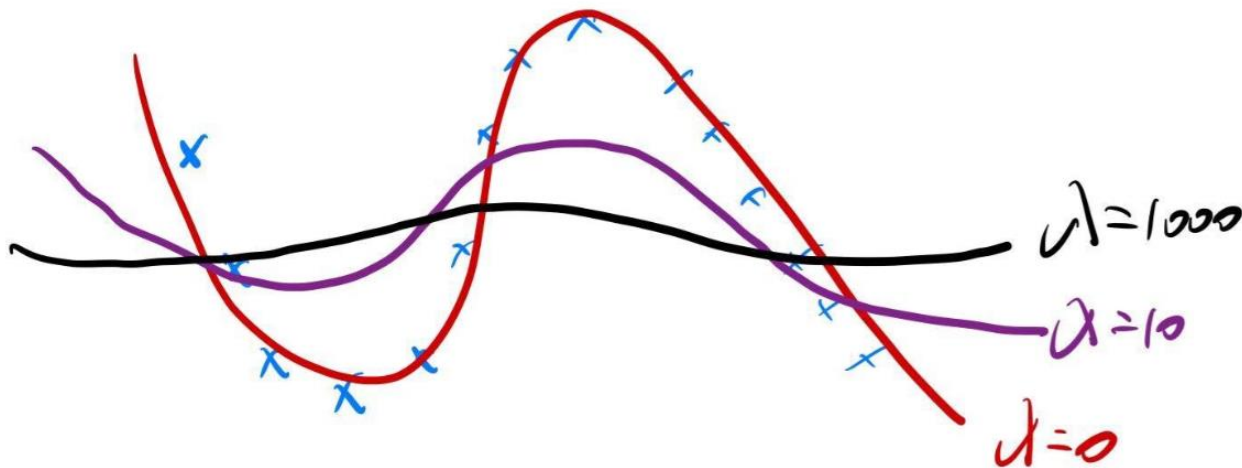
Sorry for canceling this week's tutorial due to some health issues.

- Regularization
- RBF
- Linear Classifiers

## Regularization: start from L2-norm

$$f(\mathbf{w}; \mathbf{x}) = \left\| \mathbf{X}\mathbf{w} - \mathbf{Y} \right\|_2^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

- Already explained it from “feature selection” perspective (L1, L0)
- Here we understand from the model’s capacity (L2)



$$\mathbf{y} = \mathbf{w}^T \mathbf{x}$$

$$\|\mathbf{y}\| = \|\mathbf{w}^T \mathbf{x}\| \leq \|\mathbf{w}\| * \|\mathbf{x}\|$$

# Regularization: compare L0, L1, L2

## L0- vs. L1- vs. L2-Regularization

	Sparse 'w' (Selects Features)	Speed	Unique 'w'	Coding Effort	Irrelevant Features
L0-Regularization	Yes	Slow	No	Few lines	Not Sensitive
L1-Regularization	Yes*	Fast*	No	1 line*	Not Sensitive
L2-Regularization	No	Fast	Yes	1 line	A bit sensitive

# Regularization: L2 v.s. L1

	Sparse 'w' (Selects Features)	Unique 'w'
L0-Regularization	Yes	No
L1-Regularization	Yes*	No
L2-Regularization	No	Yes

## L2-Regularization vs. L1-Regularization

- L2-Regularization:

- Insensitive to changes in data.
- Decreased variance:
  - Lower test error.
- Closed-form solution.
- Solution is unique.
- All ' $w_j$ ' tend to be non-zero.
- Can learn with *linear* number of irrelevant features.
  - E.g., only  $O(d)$  relevant features.

- L1-Regularization:

- Insensitive to changes in data.
- Decreased variance:
  - Lower test error.
- Requires iterative solver/subgradient methods.
- Solution is not unique.
- Many ' $w_j$ ' tend to be zero.
- Can learn with **exponential** number of irrelevant features.
  - E.g., only  $O(\log(d))$  relevant features.

[Paper on this result by Andrew Ng](#)

# Regularization: L1 norm

## L1-loss vs. L1-regularization

- Don't confuse the L1 loss with L1-regularization!
  - L1-loss is robust to outlier data points.
    - You can use this instead of removing outliers.
  - L1-regularization is robust to irrelevant features.
    - You can use this instead of removing features.
- And note that you can be robust to outliers and irrelevant features:

$$f(w) = \underbrace{\|Xw - y\|_1}_{L_1\text{-loss}} + \lambda \underbrace{\|w\|_1}_{L_1\text{-regularizer}}$$

- Can we smooth and use “Huber regularization”?
  - Huber regularizer is still robust to irrelevant features.
  - But it's the non-smoothness that sets weights to exactly 0.

# Regularization: L0 and L1 regularization

In linear models, setting  $w_j = 0$  is the same as removing feature 'j':

$$\hat{y}_i = w_1 x_{i1} + w_2 x_{i2} + w_3 x_{i3} + \dots + w_d x_{id}$$

$\downarrow$  set  $w_2 = 0$

$$\hat{y}_i = w_1 x_{i1} + 0 + w_3 x_{i3} + \dots + w_d x_{id}$$

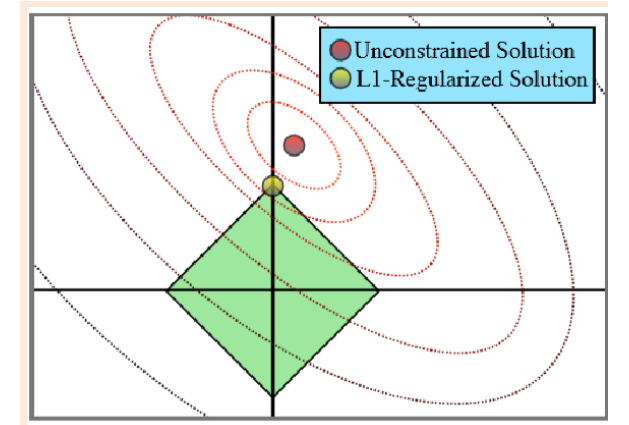
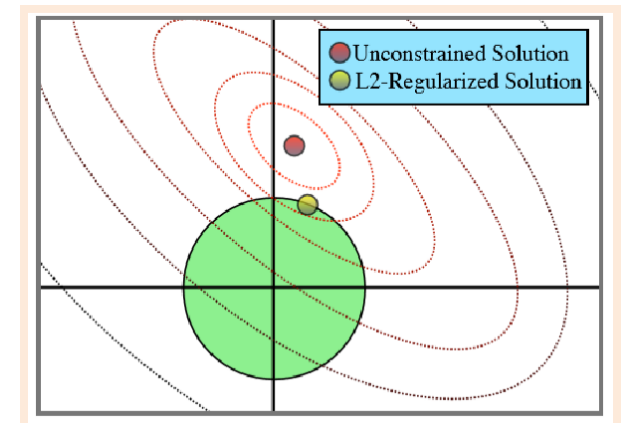
$\underbrace{\hspace{2cm}}$  ignore  $x_{i2}$

The L0 "norm" is the number of non-zero values ( $\|w\|_0 = \text{size}(S)$ ).

$$\text{If } w = \begin{bmatrix} 1 \\ 0 \\ 2 \\ 0 \\ 3 \end{bmatrix} \text{ then } \|w\|_0 = 3 \quad \text{If } w = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \text{ then } \|w\|_0 = 0.$$

- But it's **hard to find the 'w'** minimizing this objective.
- We discussed **forward selection**, but requires **fitting  $O(d^2)$  models**.
- If 'd' is large, **forward selection is too slow**:
  - For least squares, need to fit  $O(d^2)$  models at cost of  $O(nd^2 + d^3)$ .
  - **Total cost  $O(nd^4 + d^5)$** , and even if you are clever still costs  $O(nd^2 + d^4)$ .

	Speed
L0-Regularization	<b>Slow</b>
L1-Regularization	<b>Fast*</b>
L2-Regularization	<b>Fast</b>





- Regularization
- **RBF**
- Linear Classifiers

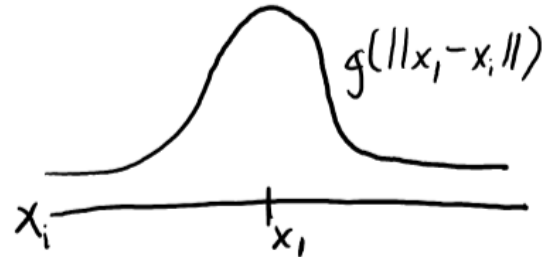
# RBF: what is Gaussian Radial Basis Functions

- What is a **radial basis functions** (RBFs)?
  - A set of non-parametric bases that **depend on distances to training points**.

Replace  $x_i = (\underbrace{x_{i1}, x_{i2}, \dots, x_{id}}_{\text{'d' features}})$  with  $z_i = (\underbrace{g(\|x_i - x_1\|), g(\|x_i - x_2\|), \dots, g(\|x_i - x_n\|)}_{\text{'h' features}})$

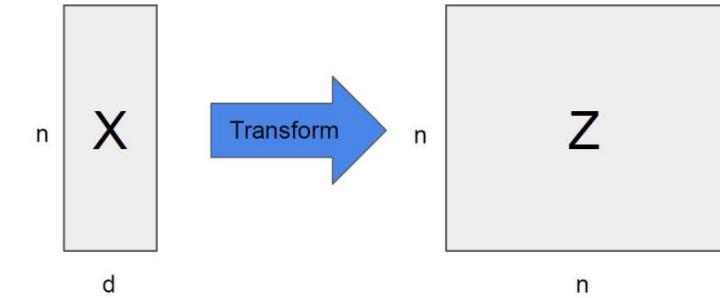
*Handwritten notes:*  
- "distance of Feature vector to example 1" points to  $g(\|x_i - x_1\|)$   
- "Distance to example 2" points to  $g(\|x_i - x_2\|)$

- Have 'n' features, with **feature 'j' depending on distance to example 'i'**.
  - Typically the **feature will decrease as the distance increases**:



# RBF: implementation

- Feature  $i$  is the “similarity” of the current example to the  $i$ -th data.



$$x_i = [x_{i1} \ x_{i2} \ \cdots \ x_{id}]$$

$$z_i = [z_{i1} \ z_{i2} \ \cdots \ z_{in}]$$

where

$$z_{ij} = g(\|x_i - x_j\|^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right)$$

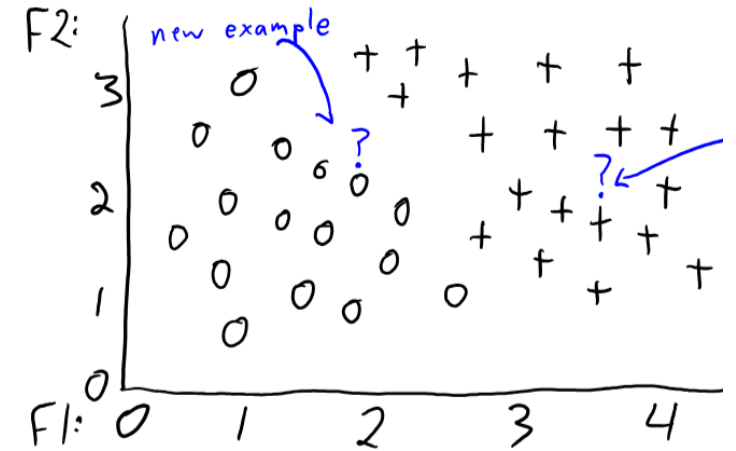
$$\tilde{X} = \left[ \begin{array}{c} \vdots \\ x_i \\ \vdots \end{array} \right] \Bigg\}^t$$

$$\tilde{Z} = \left[ \begin{array}{c} \vdots \\ g(\|\tilde{x}_i - x_j\|) \\ \vdots \end{array} \right] \Bigg\}^t$$

$\underbrace{\hspace{10em}}_n$

*Number of "features" is number of training examples*

- Also: Recall KNN



- Pseudo code:

```
z = np.zeros(n,n)
for i in range(n):
    for j in range(n):
        z[i,j] = np.exp(-norm(x[i,:], x[j,:])/2*eta^2)
```

```
z_test = np.zeros(t,n)
for i in range(t):
    for j in range(n):
        z_test[i,j] = np.exp(-norm(x_test[i,:], x[j,:])/2*eta^2)
```

# RBF: implementation

## RBFs, Regularization, and Validation

- A model that is hard to beat:
  - RBF basis with L2-regularization and cross-validation to choose  $\sigma$  and  $\lambda$ .
  - Flexible non-parametric basis, magic of regularization, and tuning for test error.

- Pseudo code:

```
for each value of  $\lambda$  and  $\theta$ :  
     $Z = \text{get\_z\_matrix}(X, \lambda, \theta)$   
     $w = (Z^T Z + \lambda I)^{-1} Z^T y$   
     $Z_{\text{val}} = \text{get\_z\_matrix}(X_{\text{val}}, \lambda, \theta)$   
     $y_{\text{vpred}} = Z_{\text{val}} * w$   
     $\text{val\_err} = \text{np.norm}(y_{\text{vpred}}, y_{\text{val}})$   
    update  $\lambda, \theta$  with the current best val_err
```

## RBF: related to linear regression

### Gaussian RBFs: A Sum of "Bumps"

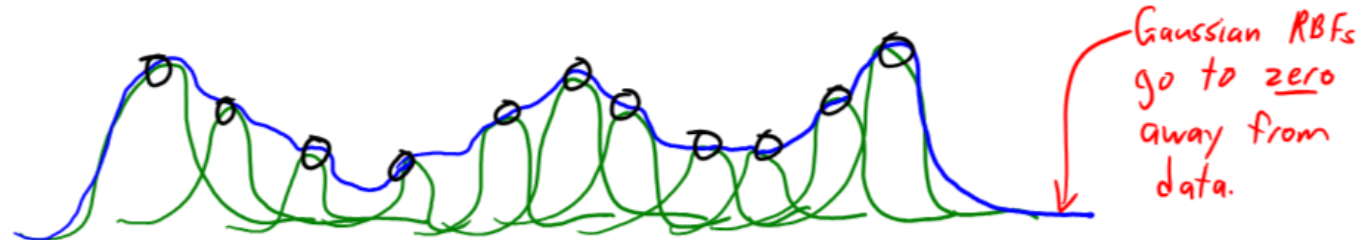
$$\hat{y}_i = w_0 \boxed{\text{flat line}} + w_1 \boxed{\text{diagonal line}} + w_2 \boxed{\text{parabola}} + w_3 \boxed{\text{cubic curve}} + w_4 \boxed{\text{U-shape}}$$

Polynomial basis represents function as sum of global polynomials.

$$\hat{y}_i = w_0 \boxed{\text{bump 1}} + w_1 \boxed{\text{bump 2}} + w_2 \boxed{\text{bump 3}} + w_3 \boxed{\text{bump 4}} + w_4 \boxed{\text{bump 5}}$$

Gaussian RBFs represent function as sum of local "bumps"

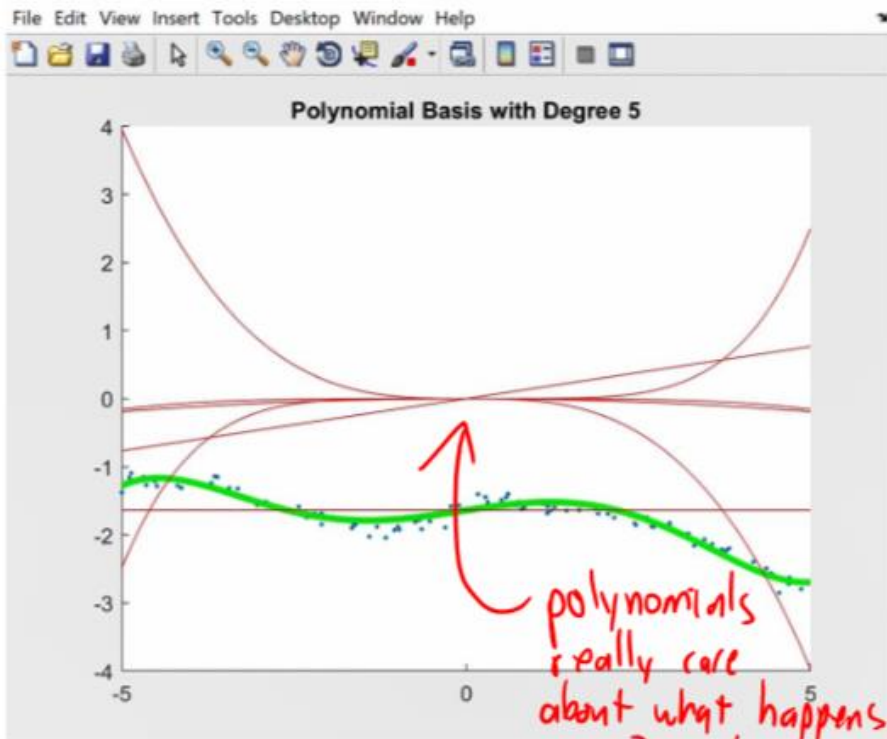
- Constructing a function from bumps ("smooth histogram"):



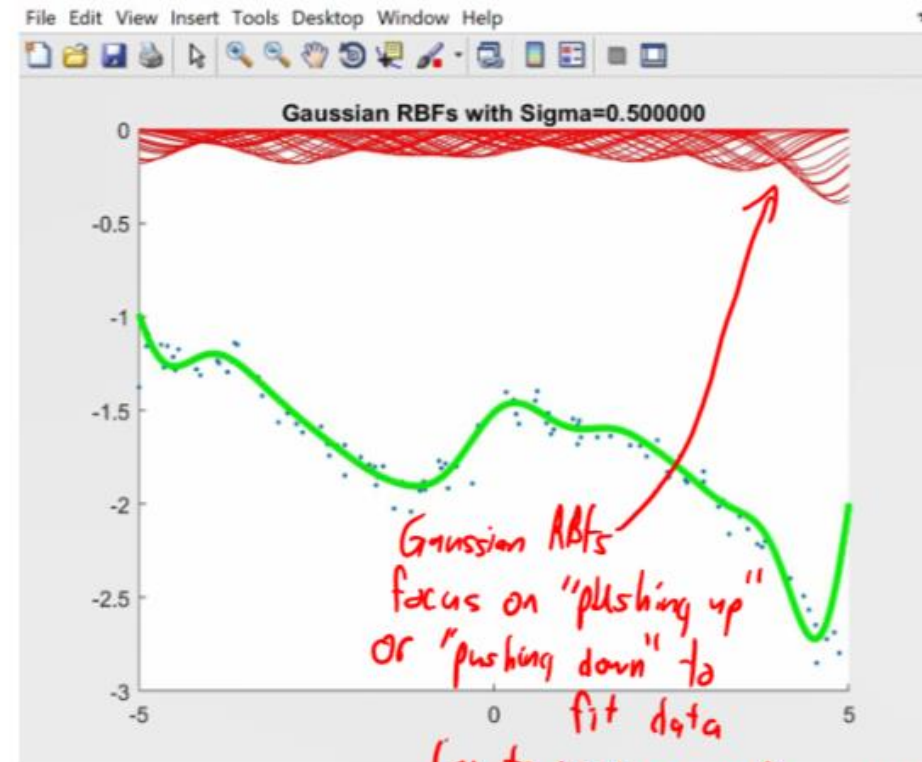
RBF: related to linear regression

## Gaussian RBFs: A Sum of “Bumps”

- Red is weight\*feature, green is prediction (sum of red lines):



polynomials really care about what happens near 0, and are wacky at edges of data

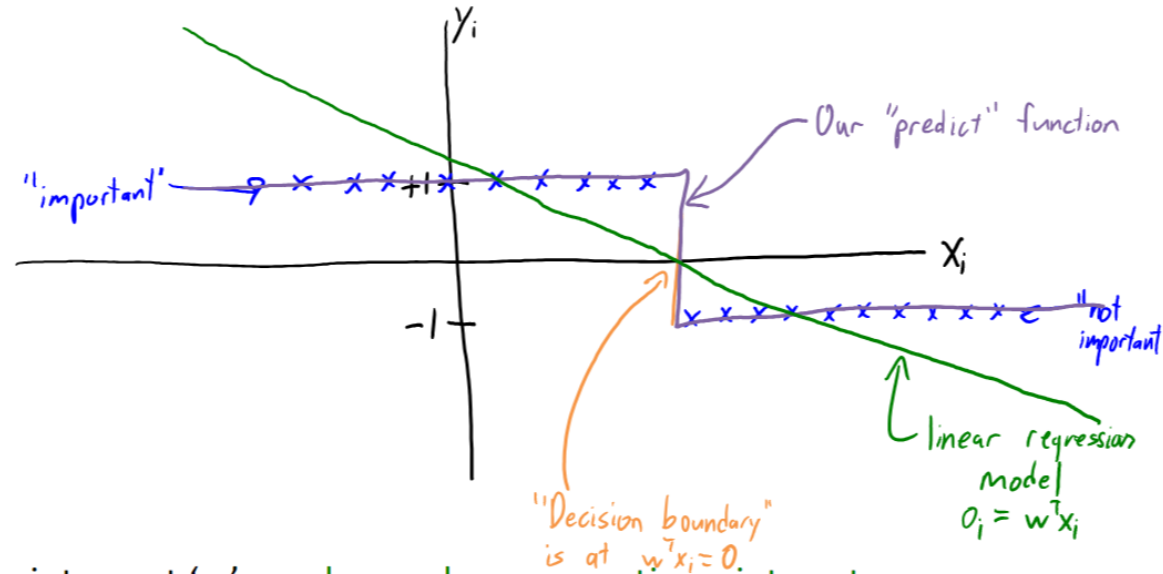


Gaussian RBFs focus on "pushing up" or "pushing down" to fit data (go to zero away from data)

- Regularization
- RBF
- **Linear Classifiers**

# Linear Classification: what is decision boundary

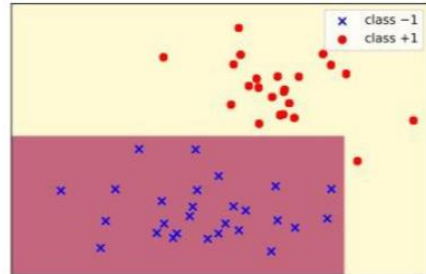
- Decision Boundary in 1D



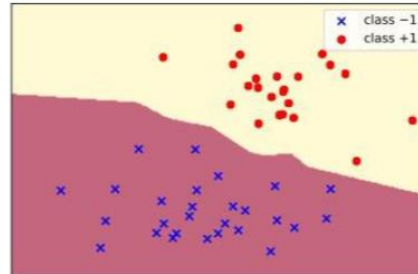
- We can interpret 'w' as a **hyperplane** separating x into sets:
  - Set where  $w^T x_i > 0$  and set where  $w^T x_i < 0$ .

- Decision Boundary in 2D

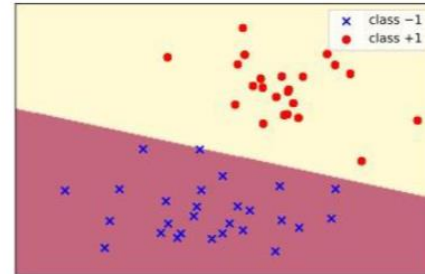
decision tree



KNN



linear classifier





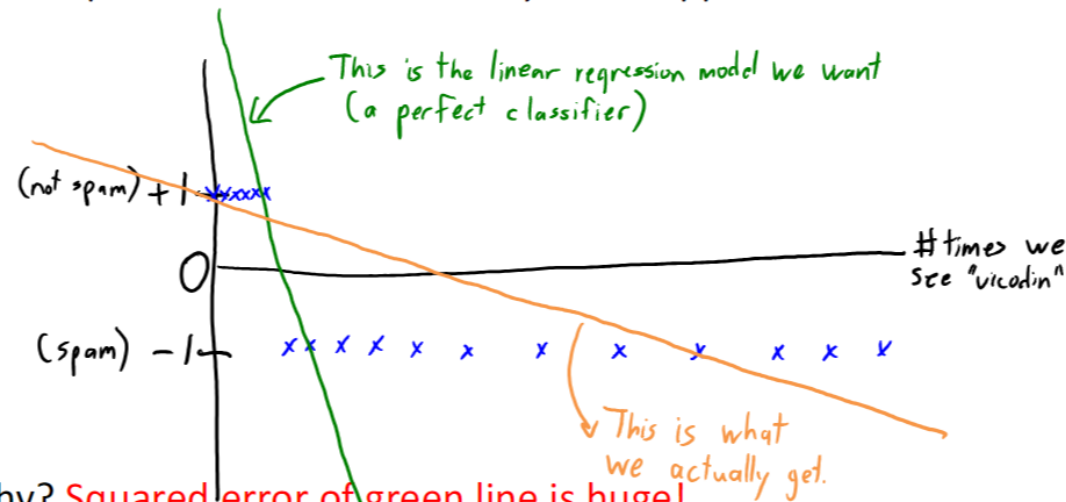
# Linear Classification: why using regression directly is not good

- Consider training by minimizing squared error with  $y_i$  that are +1 or

$$f(w) = \frac{1}{2} \|Xw - y\|^2$$

$y = \begin{bmatrix} +1 \\ +1 \\ -1 \\ +1 \\ -1 \end{bmatrix}$

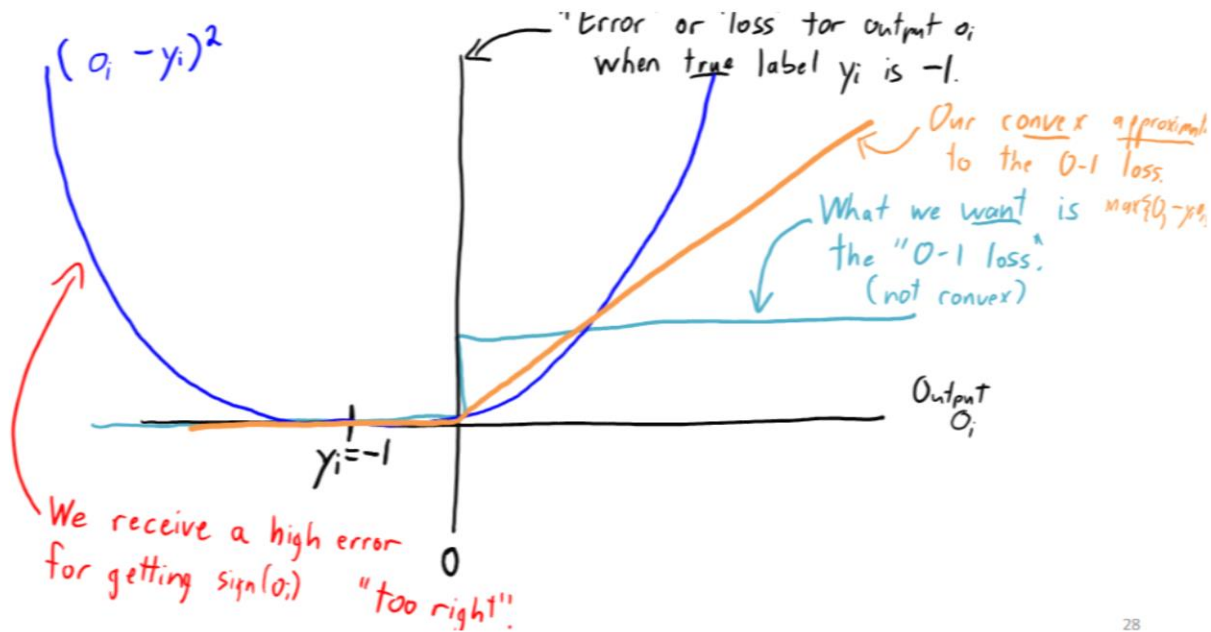
- If we predict  $o_i = +0.9$  and  $y_i = +1$ , error is small:  $(0.9 - 1)^2 = 0.01$ .
- If we predict  $o_i = -0.8$  and  $y_i = +1$ , error is bigger:  $(-0.8 - 1)^2 = 3.24$ .
- If we predict  $o_i = +100$  and  $y_i = +1$ , error is huge:  $(100 - 1)^2 = 9801$ .
  - But it shouldn't be, the prediction is correct.
- Least squares can behave weirdly when applied to classification:



- Why? Squared error of green line is huge!
  - Make sure you understand why the green line achieves an error of 0.

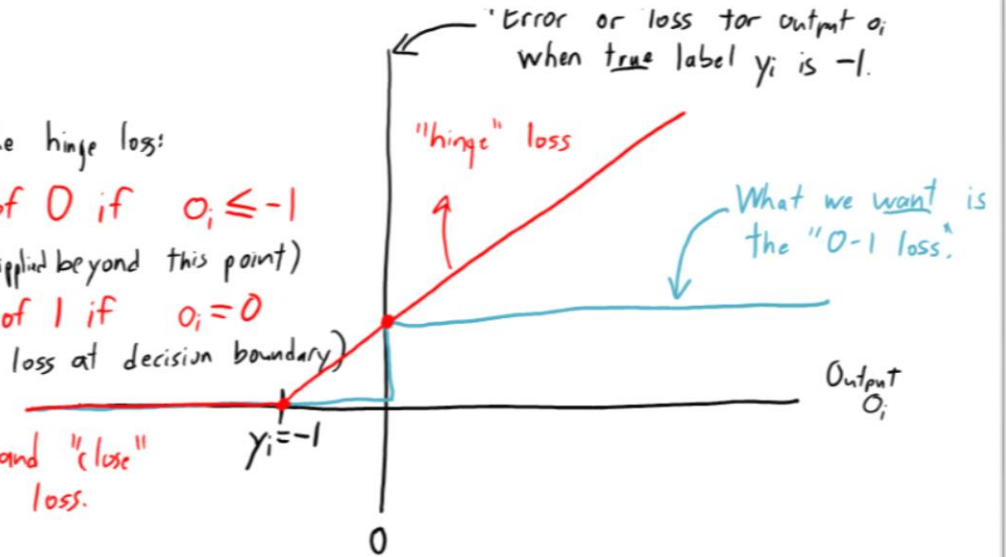
# Linear Classification: if not L2 loss, what should we use?

$$f(w) = \sum_{i=1}^n \max\{0, 1 - y_i w^T x_i\}$$



Properties of the hinge loss:

1. Has error of 0 if  $a_i \leq -1$   
(no penalty applied beyond this point)
2. Has a loss of 1 if  $a_i = 0$   
(matches 0-1 loss at decision boundary)
3. Is convex and "close" to 0-1 loss.



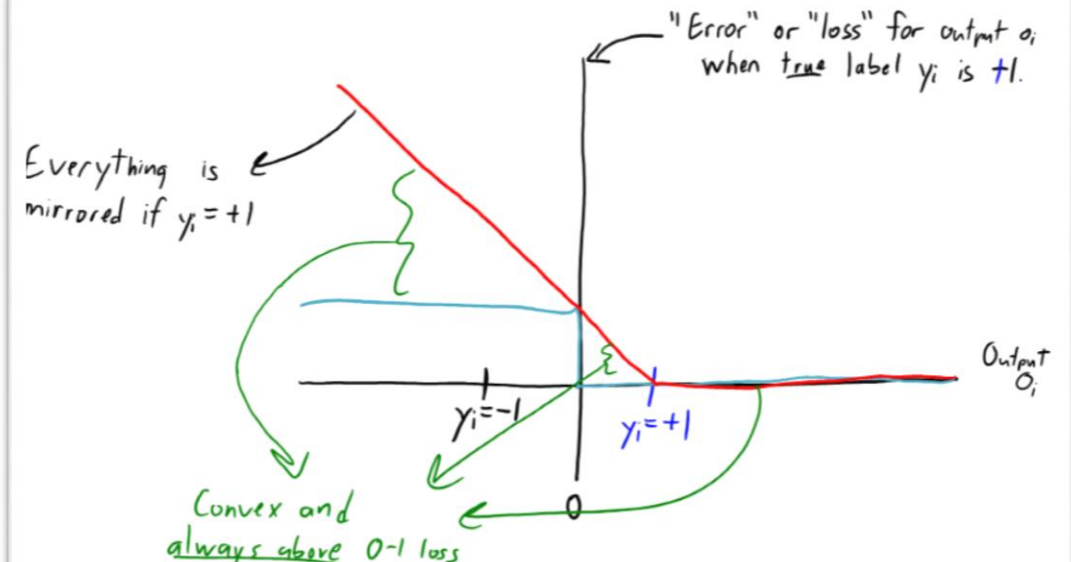
We could train by minimizing sum over all examples:

$$f(w) = \sum_{i=1}^n \max\{0, -y_i w^T x_i\}$$

But this has a **degenerate solution**:

– We have  $f(0) = 0$ , and this is the lowest possible value of 'f'.

There are two standard fixes: **hinge loss** and **logistic loss**.



Thanks for your time!

Questions?