

Linear Regression

Program:

$$\arg \min_w \frac{1}{2} \| \underset{\mathbb{R}^N}{Y} - \underset{\mathbb{R}^{N \times d+1}}{X^T w} \|_2^2$$

Solution:

$$w^* = (X^T X)^{-1} X^T Y$$

with Regularization:

Prog: $\arg \min_w \frac{1}{2} \| Y - X^T w \|_2^2 + \frac{C}{2} \| w \|_2^2$

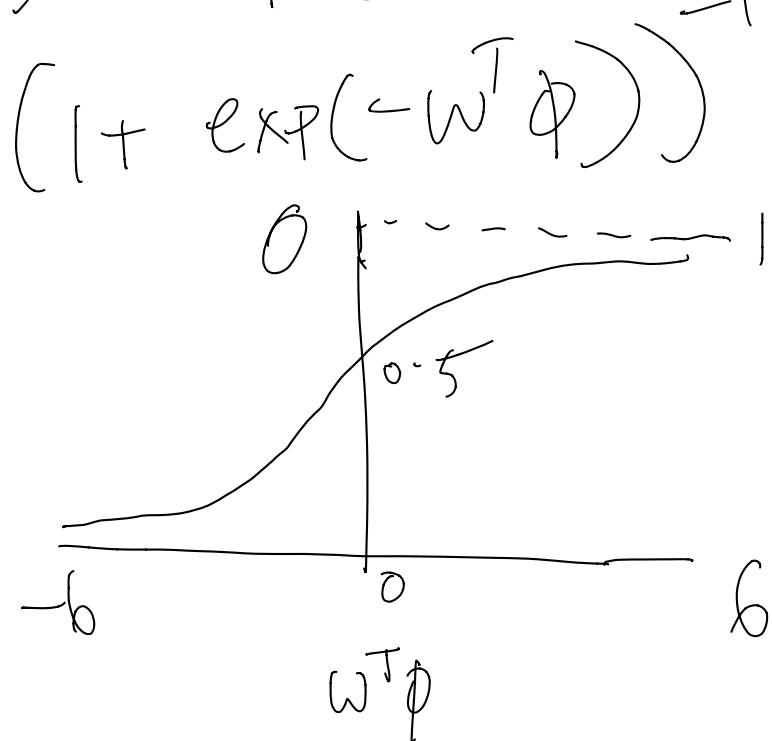
Solution: $w^* = (X^T X + CI)^{-1} X^T Y$

L1: Lasso \rightarrow Shrink features to 0
 (feature Selection)

L2: Ridge \rightarrow Shrink features but
 doesn't remove them

Logistic Regression

- Used in Classification
- Does NOT penalize when samples are easy to classify
- Sigmoid Function



- Squishes output to be between 0 and 1
- Model:

$$y^{(i)} \in \{-1, 1\}$$

$$P_{\text{err}}(i) \rightarrow$$

$$P(y^{(i)}=1) = \frac{1}{1 + \exp(-\omega^T \phi)}$$

$$P(y^{(i)}=-1) = 1 - P(y^{(i)}=1) = \frac{1}{1 + \exp(\omega^T \phi)}$$

$$P(y^{(i)} | x^{(i)}) =$$

$$\frac{1}{1 + \exp(-y^{(i)} \omega^T \phi)}$$

~ Maximum Likelihood:

$$\arg \max_w \prod P(y^{(i)} | x^{(i)})$$

$$(x^{(i)}, y^{(i)}) \in D$$

$$= \arg \min_w - \sum P(y^{(i)} | x^{(i)})$$

$$= \arg \min_w \sum_{|D|} -\log(P(y^{(i)} | x^{(i)}))$$

$$= \arg \min_w \sum_{|D|} \ln(1 + \exp(-y^{(i)} \omega^T \phi))$$

$$w \leftarrow w + \eta g(w) \quad \text{for } j \in \mathcal{I}$$

$|D|$

- Optimize Via Gradient Descent

$$\nabla_w = \sum_{|D|} \frac{-y^{(i)} \exp(-y^{(i)} w^T \phi)}{1 + \exp(-y^{(i)} w^T \phi)} \phi$$

Initialize $t=0$, w_t and α

Iterate:

$$g_t = \nabla_w f(w_t)$$

$$w_{t+1} = w_t - \alpha g_t$$

$$t \leftarrow t + 1$$

Optimization

When Can we find optimum?

- Linear, Least Squares, Convex etc
- Generally can be hard
- Jensen's Inequality:

Convex:

$$f((1-\lambda)w_1 + \lambda w_2) \leq (1-\lambda)f(w_1) + \lambda f(w_2)$$

- $\nabla^2 f(w) \succeq 0 \forall w$
- Batch vs. Stochastic
 - Batch looks at whole data set
 - Stochastic looks at a part of the data set.

$$\min_w f_0(w)$$

$$\text{S.t. } f_i(\omega) \leq 0 \quad i \in \{1, \dots, c_1\}$$

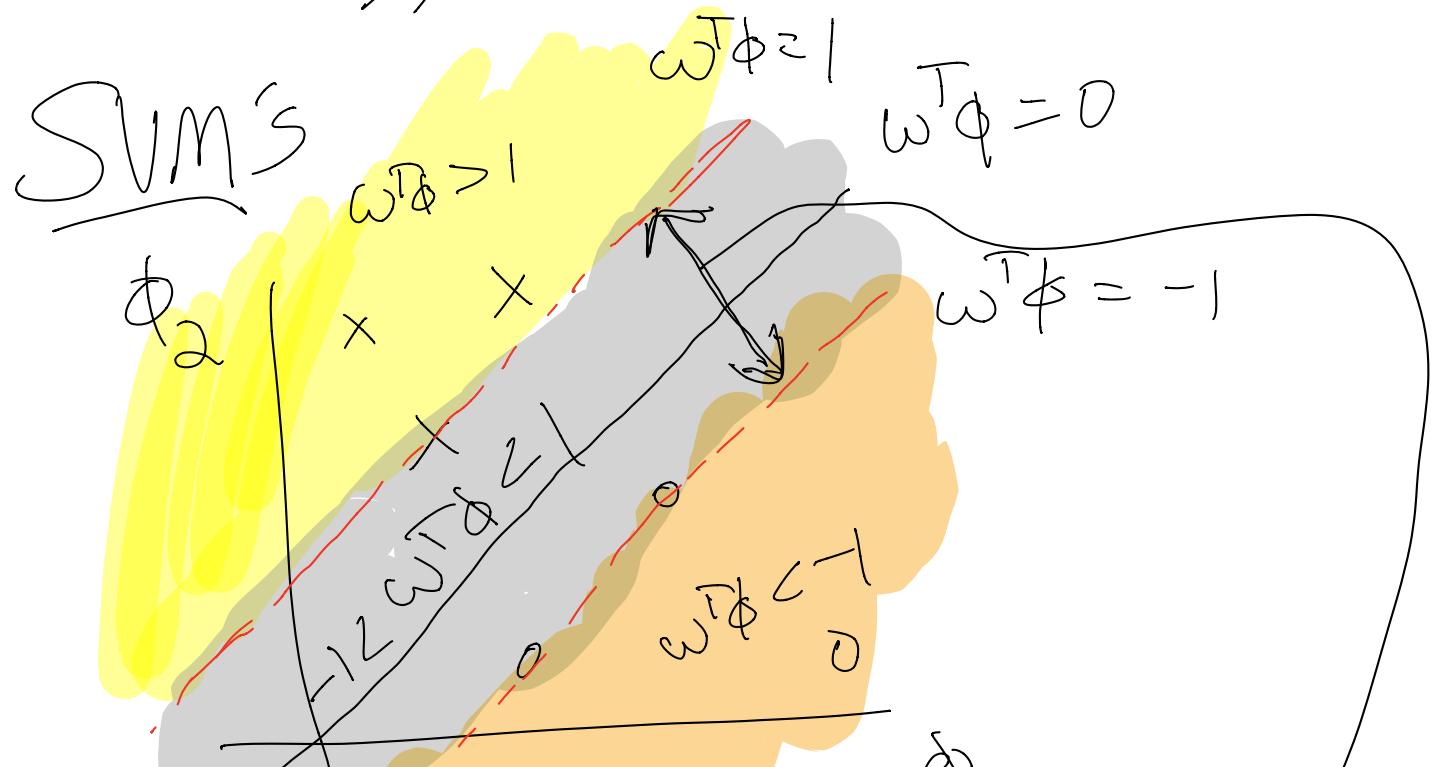
$$h_i(\omega) = 0 \quad i \in \{1, \dots, c_2\}$$

$$L(\omega, \gamma, \nu) = f_0(\omega) + \sum_{i=1}^{c_1} \gamma_i f_i(\omega) + \sum_{i=1}^{c_2} \nu_i h_i(\omega)$$

≥ 0

$$f(\omega) \geq L \geq \min_{\omega} L = g(\gamma, \nu)$$

$$\max_{\gamma, \nu} g(\gamma, \nu) \text{ s.t. } \gamma_i \geq 0 \quad \forall i$$



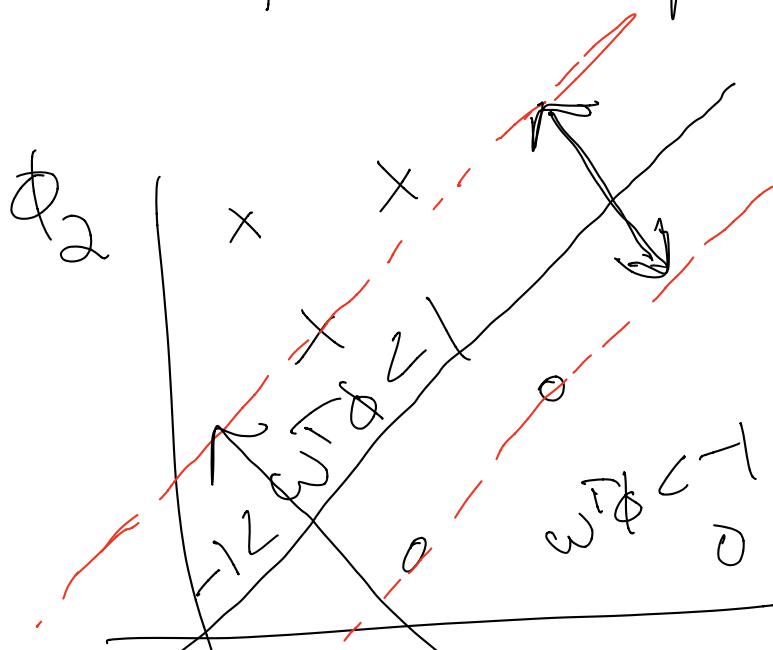
goal of SVM is to max this margin
 (separate data the most).

$$\text{margin} = \frac{2}{\|w\|}$$

$$\max_w \frac{2}{\|w\|} = \underbrace{\min_w \frac{1}{2} \|w\|_2^2 \text{ ST } y_i w^T \phi \geq 1}$$

• Slack Variable (ϵ)

- for linearly unseparable data.



$$\min_{\omega, \varepsilon \geq 0} \frac{C}{2} \|\omega\|_2^2 + \sum_{i \in D} \varepsilon^{(i)}$$

ST. $y^{(i)} \omega^T \phi \geq 1 - \varepsilon^{(i)}$

$$\varepsilon \geq 1 - y^i \omega^T \phi, \varepsilon \geq 0$$

$$\min_{\omega} \frac{C}{2} \|\omega\|_2^2 + \sum \max(0, 1 - y^{(i)} \omega^T \phi)$$

takes care
of $\varepsilon \geq 0$

$$L(\cdot) = \frac{C}{2} \|\omega\|_2^2 + \sum_i \varepsilon^i + \sum_i \alpha^i (1 - y^i \omega^T \phi - \varepsilon^i)$$

$$= \frac{C}{2} \|\omega\|_2^2 - \omega^T \sum_i \alpha^i y^i \phi + \sum_i \varepsilon^i (1 - \alpha^i) + \sum_i \alpha^i$$

$$\frac{\partial L_a}{\partial \omega^0} = C \omega = \sum_i \alpha^i y^i \phi$$

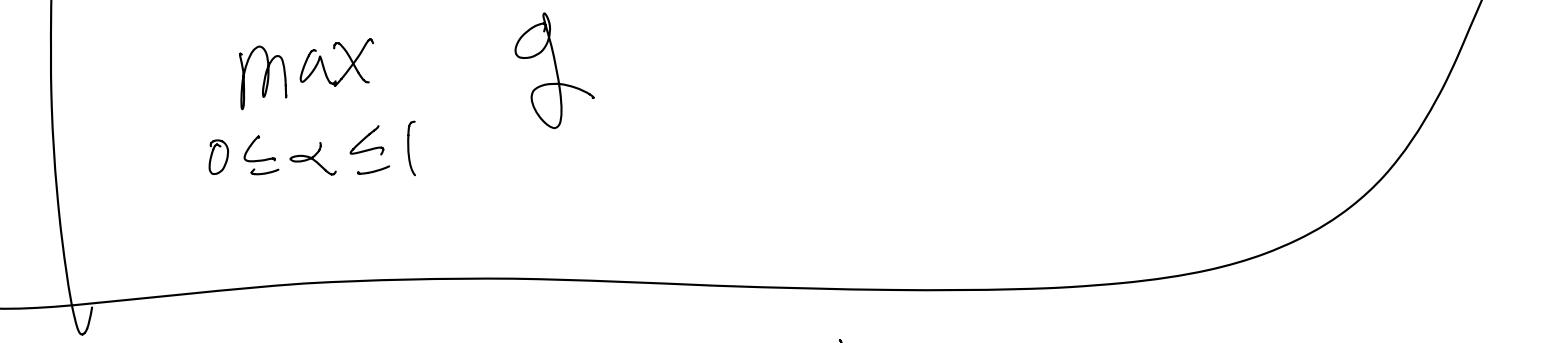
$$\frac{\partial L}{\partial \alpha} = \min_{\alpha \geq 0} \sum_i (\alpha - \alpha^i) \rightarrow \alpha \leq 1$$

b/c otherwise
 $\alpha \rightarrow \infty$ as this
 would make
 $\sum_i (\alpha - \alpha^i) \rightarrow \infty$

$$g = \frac{c}{2} \left(\frac{1}{C} \left\| \sum_i \alpha^i y^i \phi \right\|_2^2 - \frac{1}{C} \sum_i \alpha^i y^i \phi \sum_i \alpha^i y^i \phi \right)$$

$$\begin{aligned}
 & + \sum_i \alpha^i \\
 & = \cancel{\frac{c}{2}} \left(\frac{1}{C} \left\| \sum_i \alpha^i y^i \phi \right\|_2^2 \right. \\
 & \quad \left. - \frac{1}{C} \left\| \sum_i \alpha^i y^i \phi \right\|_2^2 + \sum_i \alpha^i \right)
 \end{aligned}$$

$$= -\frac{1}{2} C \left\| \sum_i \alpha^i y^i \phi \right\|_2^2 + \sum_i \alpha^i$$

$$\max_{0 \leq \alpha \leq 1} g$$


Log loss is Hinge loss:

$$\lim_{\epsilon \rightarrow 0} \epsilon \log \left(1 + \exp \frac{-F}{\epsilon} \right)$$

$$\text{when } F \geq 0 = 0$$

$$F \geq 0$$

$$\text{when } F \leq 0 = -F$$

$$F \leq 0$$

$$\lim_{\epsilon \rightarrow 0} \epsilon \log \left(1 + \exp \frac{-F}{\epsilon} \right)$$

$$= \max(0, -F)$$

SVM is 0-temp limit of logistic regression

$$\min_{\omega} \frac{1}{2} \|\omega\|_2^2 + \sum c \log \left(1 + \exp \left(\frac{L - y \omega \phi}{c} \right) \right)$$

$c \rightarrow 0 \rightarrow$ SVM

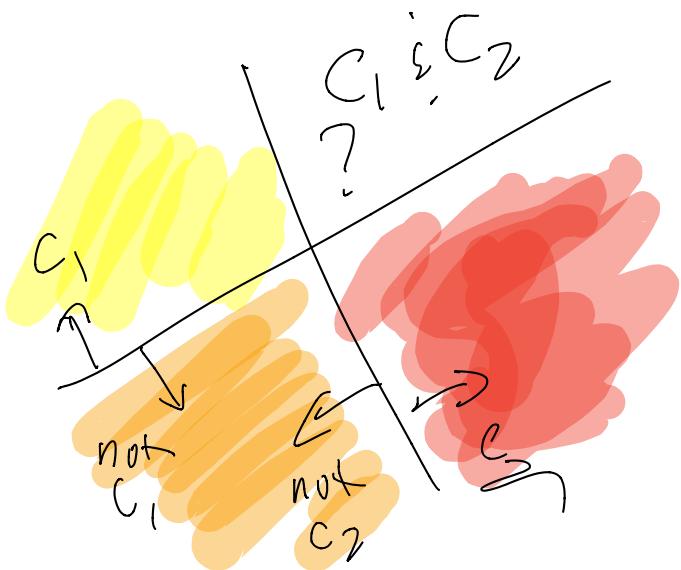
$c = 1 \rightarrow$ Log Reg.

Multi Class

- 1 vs All

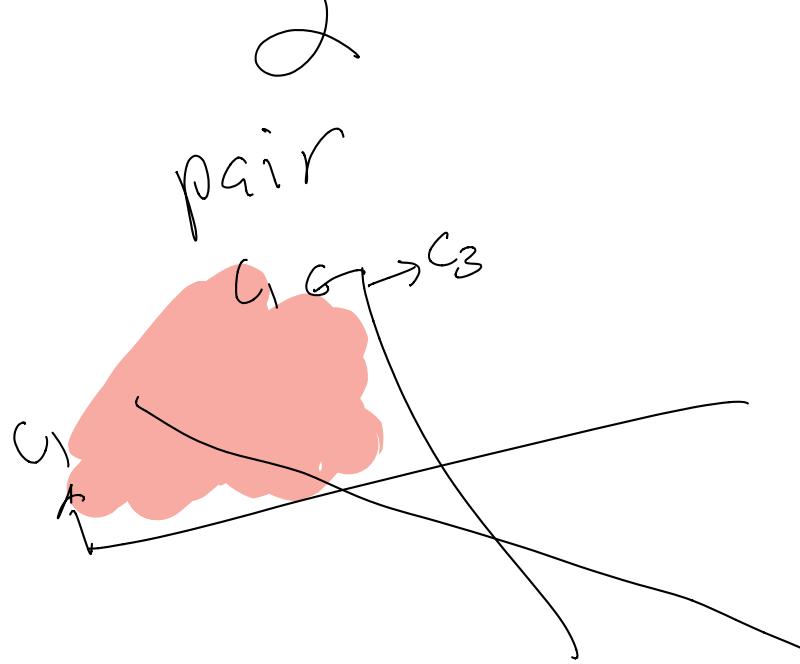
- K-1 classifiers
- is point in a class or

Not



- 1 vs 1

- $\underbrace{K(K-1)}$, one for each



Classify w/ majority vote

— Multi class Logistic Regression

$$P(y^{(i)} = k | x^{(i)}) = \frac{\exp w_{(k)}^T \phi(x^{(i)})}{\sum_{j \in \{0, \dots, K-1\}} \exp w_{(j)}^T \phi(x^{(i)})}$$

interpret as probabilities

$$\arg \min_w \sum -\log P(y=y^{(i)} | x^{(i)})$$

Minimize the sum of $-\log$ of the correct class Probability.

- Multi Class SVM

$$w_{y^{(i)}}^T \phi \geq w_y^T \phi \quad \text{for all } y \in \{0, \dots, K-1\}$$

- we want correct class score ($w_{y^{(i)}}^T \phi$) to be bigger than all other scores for every other class

$$\min_w \frac{1}{2} \|w\|_2^2 + \sum_i \max_{\hat{y}} \left(1 + w^T \varphi(x^i, \hat{y}) \right) - w^T \varphi(x^i, y^i)$$

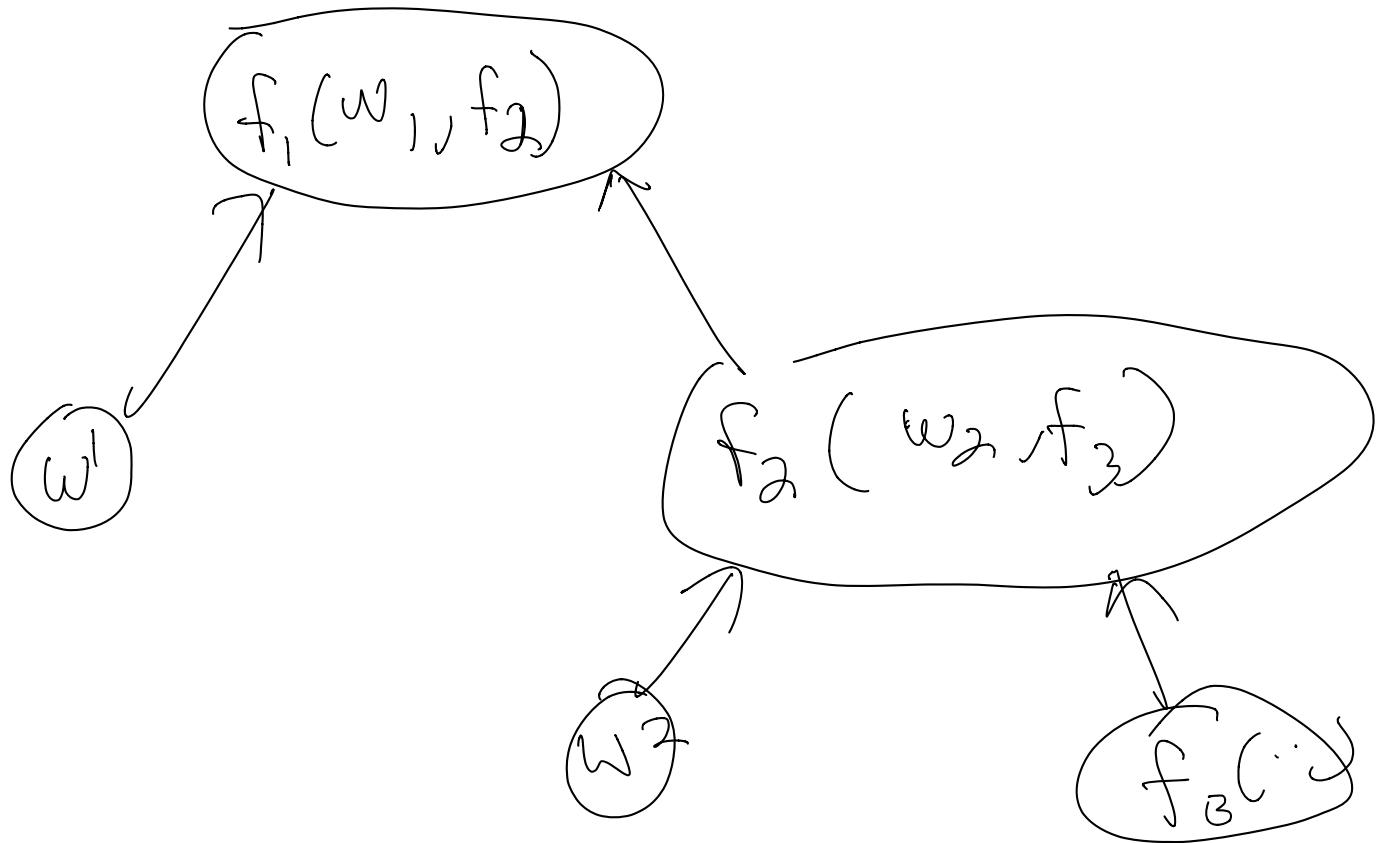
Score our model gives (inference).

Score of correct class

- Generally

$$\min_w \frac{1}{2} \|w\|_2^2 + \sum_i \ln \sum_{\hat{y}} \exp \frac{L(y^i, \hat{y}) + w^T \varphi(x^i, \hat{y})}{e} = w^T \varphi(x^i, y^i)$$

Deep Nets



- ~ Conv Nets
 - decrease Spatial resolution and # channels.
- ~ Fully Connected
 - $wX + b$
 - Trainable P's: $w \in b$
 - Issues b/c if $X \in 256 \times 256$, you will have ALOT of weights to tune.

- Convolutional Layer

- decrease spatial resolution and # channels.



$$W_2 = \frac{W_1 - F + 2P}{S} + 1$$

$K = \# \text{ filters}$
 $F = \text{size}$
 $S = \text{stride}$
 $P = \text{padding}$

$$h_2 = \frac{h_1 - F + 2P}{S} + 1$$

$$D_2 = K$$

$F \cdot F \cdot D_1$ weights/filter

$F \cdot F \cdot D_1 \cdot K$ weights total

- Max Pooling

- reduce spatial size
- prevent overfitting

- runs a window along an image and replaces the values in the window w/ max / average



$$\omega_2 = \frac{\omega_1 - F}{S} + 1$$

$$H_2 = \frac{H_1 - F}{S} + 1$$

— Drop Out

- randomly set activations to 0
- regularization

— Use NLL loss w/ softmax

— Train using Stochastic Gradient Descent e) momentum-

$$F(w, x, Y) =$$

$$f_1(w_1, y, f_2(w_2, f_3(w_3, x)))$$

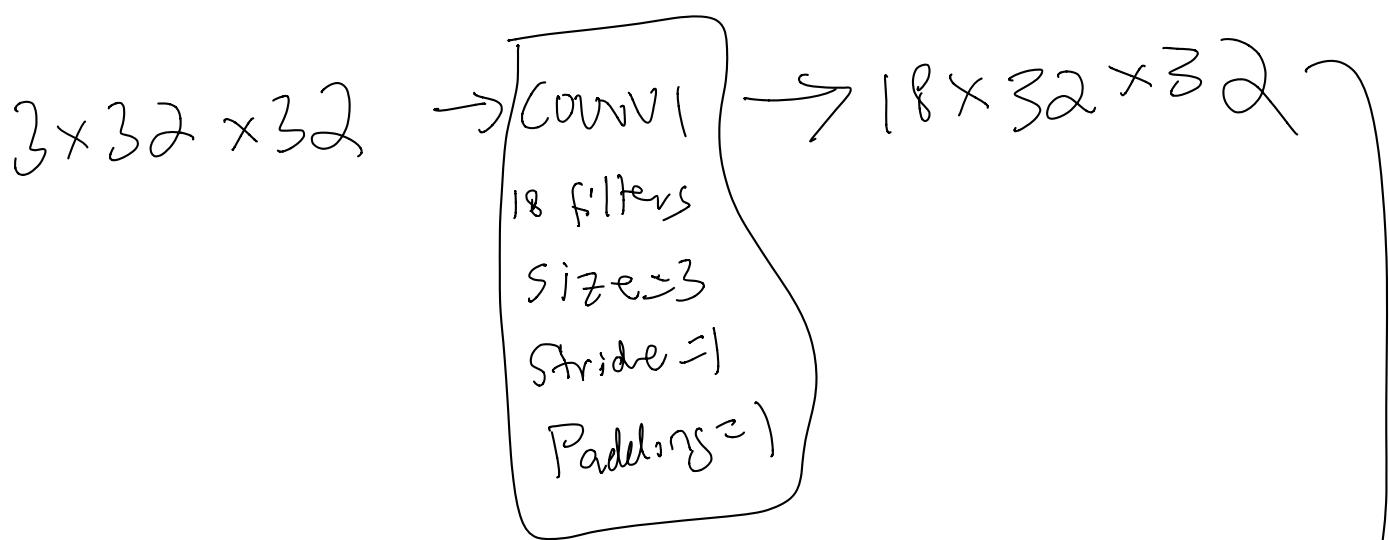
$$x_3 = f_3(w_3, x)$$

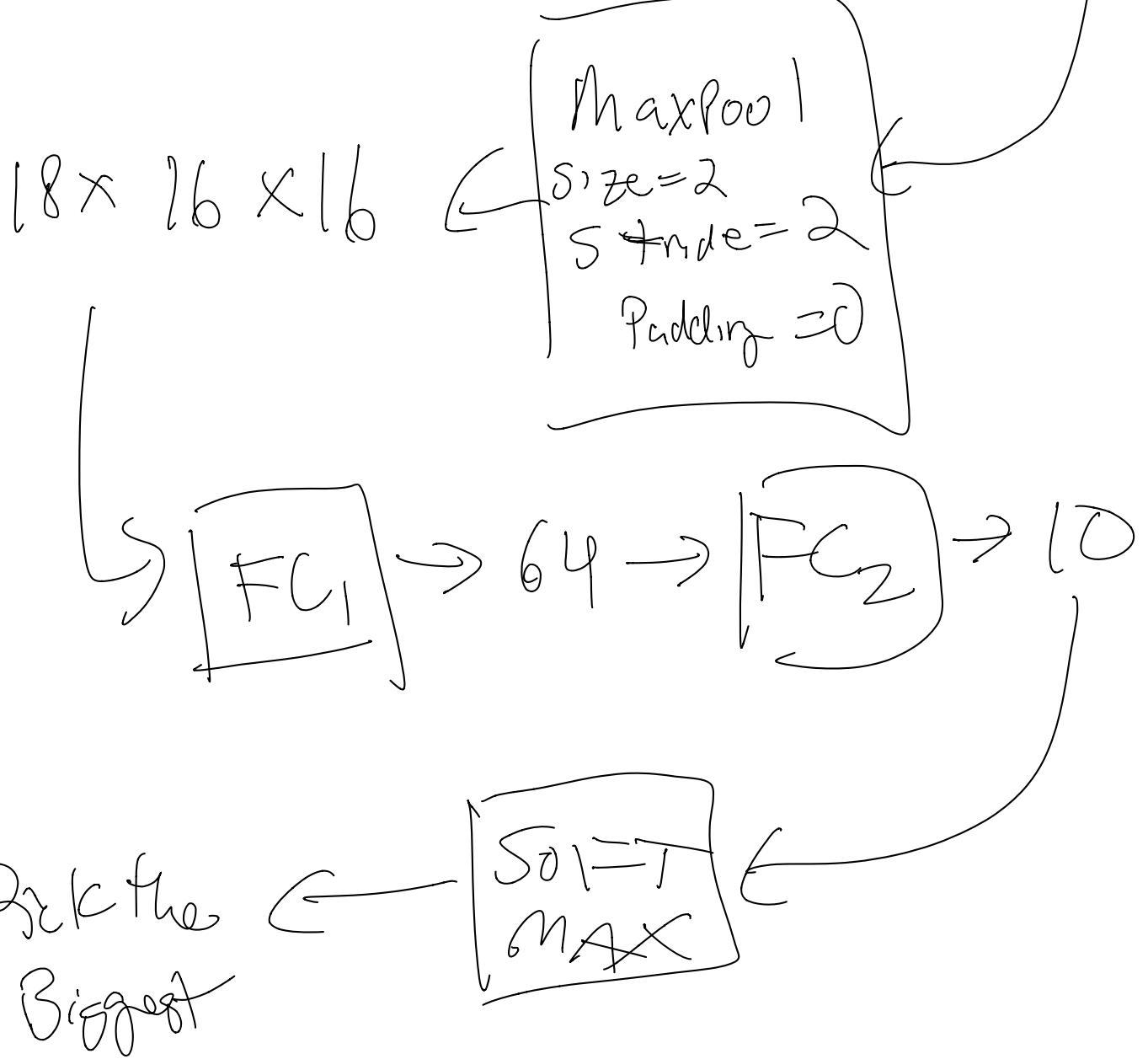
$$x_2 = f_2(w_2, x_2)$$

$$\frac{\delta F}{\delta w_3} = \frac{\delta f_1}{\delta x_1} \cdot \frac{\delta x_1}{\delta x_2} \cdot \frac{\delta x_2}{\delta w_3}$$

$$= \frac{\delta f_1}{\delta f_2} \cdot \frac{\delta f_2}{\delta f_3} \cdot \frac{\delta f_3}{\delta w_3}$$

$$\frac{\delta F}{\delta w_2} = \frac{\delta f_1}{\delta x_1} \cdot \frac{\delta x_1}{\delta w_2} = \frac{\delta f_1}{\delta f_2} \cdot \frac{\delta f_2}{\delta w_2}$$





Boosting

- Weak classifiers working together to be a strong classifier
- Given:

$$D = \{ (x^{(i)}, y^{(i)}) \}$$

$F = \{f_t\}$ Set of weak classifiers.

$$F_T(x) = \sum_{t=1}^T \alpha_t f_t(x)$$

$$F_t = F_{t-1}(x) + \alpha_t f_t(x)$$

Decision Tree

If ensemble

- Start at the root and follow decision
- leaf node reveals results

Learning

- choose a variable that "best" splits the data

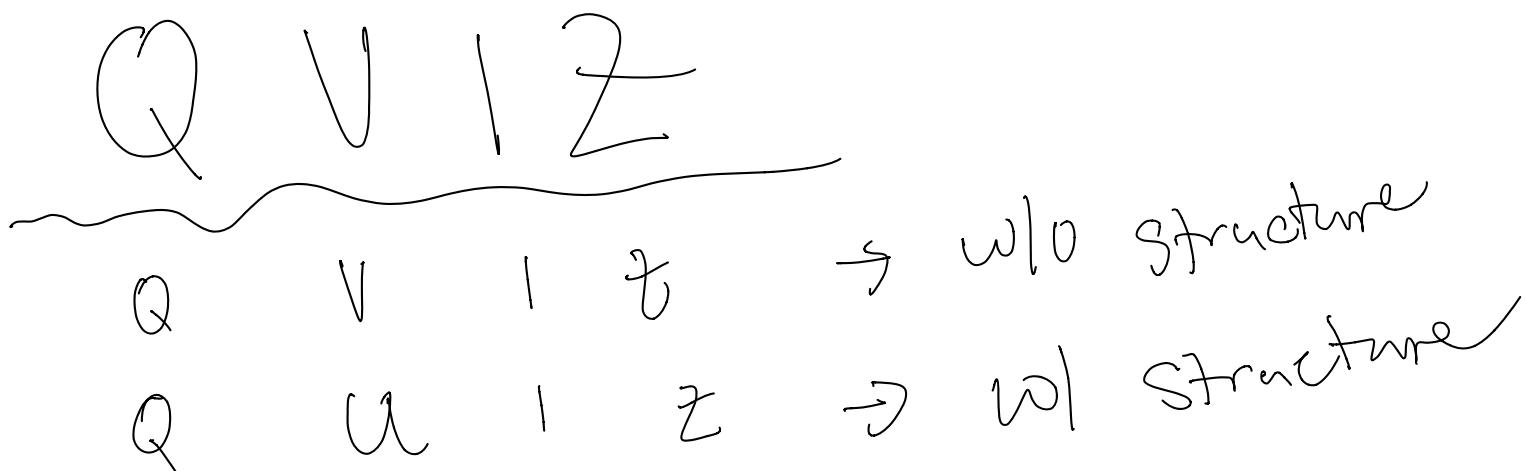
' Split data according to rule

' Stop once # is small

Structured Prediction

Inference

- find the highest score from out put
- easy for binary - 1000 classes



- Formulate it as predictions of all 4 letter words.

$$\underbrace{2^6}_{\text{#feat}} \quad \underbrace{2^6}_{\text{#feat}} \quad \underbrace{2^6}_{\text{#feat}} \quad \underbrace{2^6}_{\text{#feat}} \rightarrow 2^{\underline{6}^4}$$

- Size of Feature Space \approx # classes

- Formally

$$\mathbf{Y} = (y_1, \dots, y_D) \quad y_d \in \{1, \dots, K\}$$

$$\mathbf{Y}^* = \arg \max_{\mathbf{\hat{y}}} F(\mathbf{w}, \mathbf{x}, \mathbf{\hat{y}})$$

$$= \arg \max_{\mathbf{\hat{y}}} F(\mathbf{w}, \mathbf{x}, \hat{y}_1, \dots, \hat{y}_D)$$

Basically, pick configuration $w/$
highest score