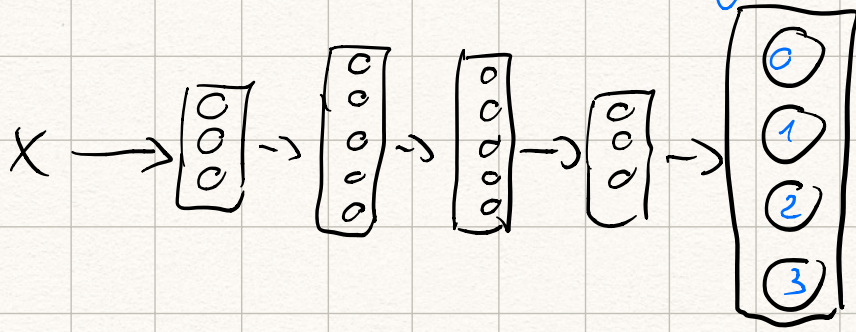


Softmax layer



$$\frac{e^5}{176.3} = 0.842$$

$$\frac{e^2}{176.3} = 0.042$$

$$\frac{e^{-1}}{176.3} = 0.002$$

$$\frac{e^3}{176.3} = 0.114$$

$$z^l = w^l a^{l-1} + b^l$$

Activation:

$$t = e^{z^l}$$

$$a^l = \frac{e^{z^l}}{\sum_{i=1}^4 t_i}$$

of classes

$$a_{ij} = \frac{t_{ij}}{\sum_{i=1}^4 t_{ij}}$$

$$a^l = g^l(z^l)$$

softmax

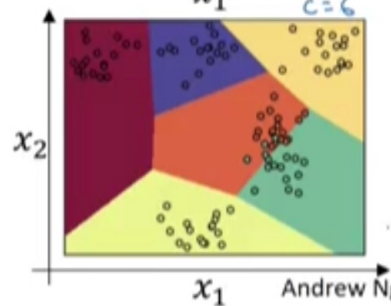
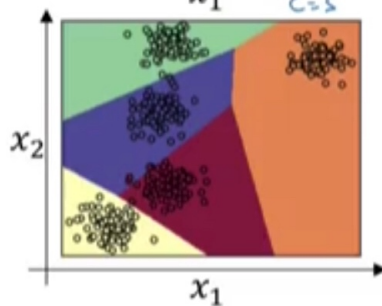
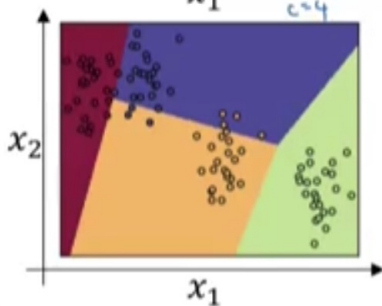
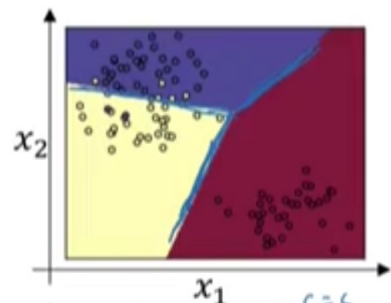
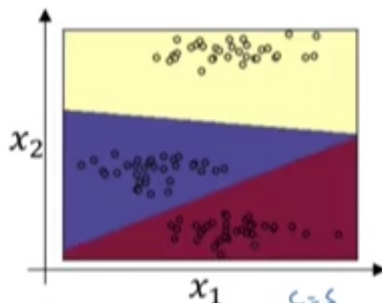
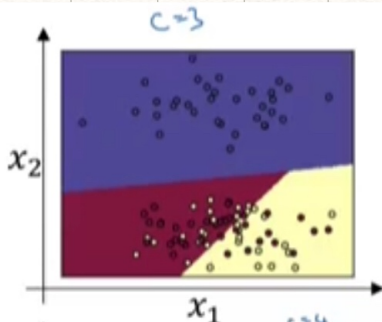
$$z^l = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}$$

$$t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix}$$

$$\Rightarrow \sum_{i=1}^4 t_i = 176.3$$

$$a^l = \frac{t}{176.3}$$

Softmax examples



Training a softmax classifier

Softmax generalizes logistic regression to C classes,

loss function:

$$L(\bar{y}, y) = - \sum_{j=1}^C y_j \log \bar{y}_j$$

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad \bar{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

$$-y_2 \log \bar{y}_2 = -\log \bar{y}_2$$

\Rightarrow make \bar{y}_2 as big as possible

ML strategy

- Fit training set well on cost func. \Rightarrow • Bigger network
• Better opt. algorithm
- Fit dev set well on cost func. \Rightarrow • Regularization
• Bigger training set
- Fit test set well on cost func. \Rightarrow • Bigger dev set
- Performs well in real world \Rightarrow change dev set or cost func.

\Rightarrow Set up a single number evaluation metric

- Precision: examples recognised as cat, what % actually are cats?
- Recall: what % of actual cats are correctly recognized?

• F1 score: $\frac{2}{\frac{1}{P} + \frac{1}{R}}$ \Rightarrow "harmonic mean"

\Rightarrow satisfying and optimizing metric

classifier	Accuracy \swarrow optimizing metric	Running time \swarrow satisfying metric
A	90%	80 ms
B	92%	95 ms
C	95%	1500 ms

\Rightarrow maximize accuracy

subject to running time ≤ 100 ms

Train/dev/test set distributions

Regions:

- US
- UK
- Europe
- India
- China
- Asia

} Dev

} Test

\Rightarrow different distribution
very bad idea

\Rightarrow randomly shuffle into dev/test set

\Rightarrow same distribution

Size of dev/test set

Old way :

$\sim 100 / 1000 / 10000$
examples

60%	20%	20%
Train	dev	test

Big data :

$\sim 10^6 \ll$ examples

98%	1%	1%
Train	dev	test

change dev/test set or metric

Metric : classification error

A : 3% error \rightarrow but shows Paragrophic as well

B : 5% error

Metric + Dev \Rightarrow prefer A

You \Rightarrow prefer B

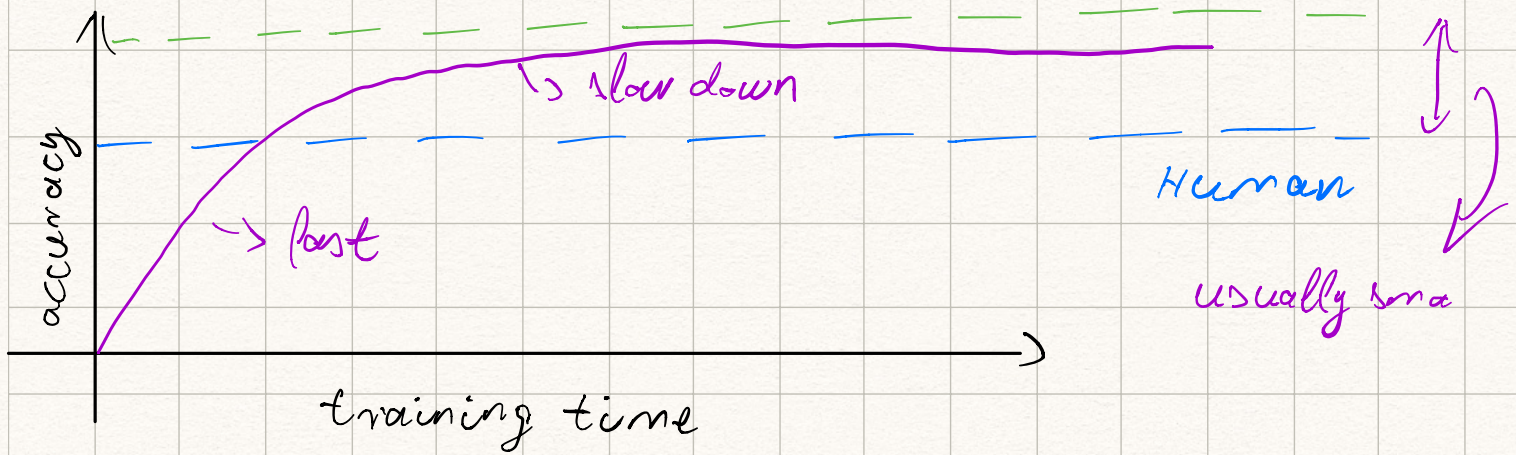
Weighting your metric

$$\text{Error} : \frac{1}{\sum_i w^i} \quad \cancel{\frac{1}{n_{\text{dev}}}} \quad \sum_{i=1}^{n_{\text{dev}}} w^i h(y_{\text{pred}}^i \neq y^i)$$

$$\Rightarrow w^i = \begin{cases} 1 & \text{if } x^i \text{ is non-porn} \\ 10 & \text{if } x^i \text{ is porn} \end{cases}$$

Human level performance

Bayes optimal error



Avoidable bias

Human	1%	↑	7.5%	↑ avoidable bias
Training error	8%	↓	8%	↓ bias
Dev error	10%		10%	

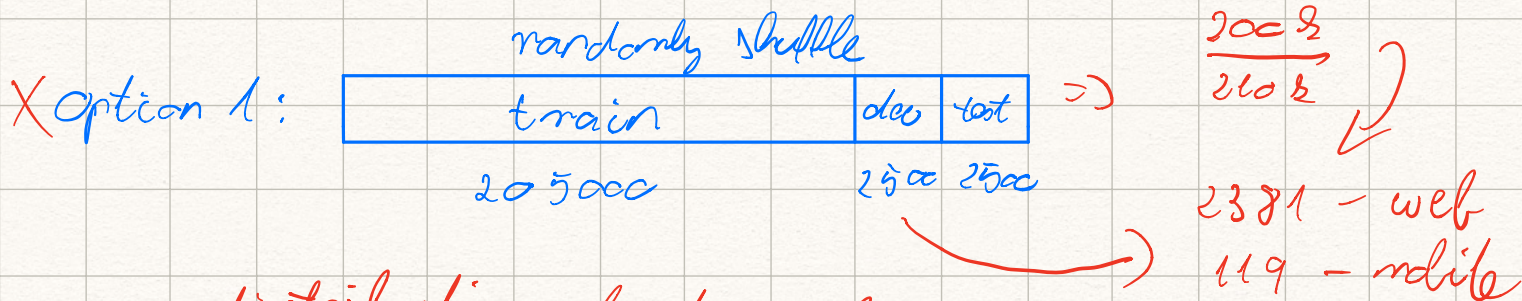
Focus on bias
Bigger networks,
better opt.
train longer

Focus on variance
Regularization →
more data
NN architecture

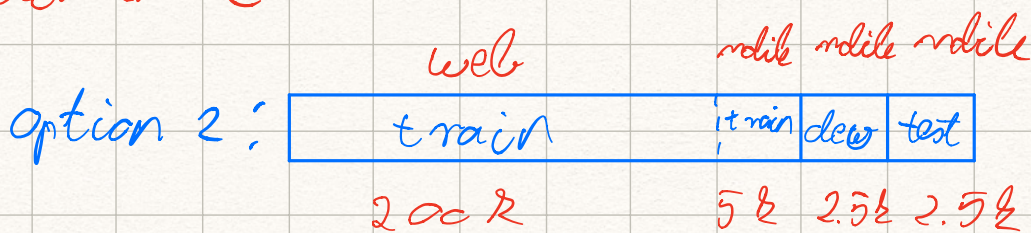
Train and test on different distribution

Data from web ~ 200000 , good quality

Data from mobile ~ 10000 , lower quality

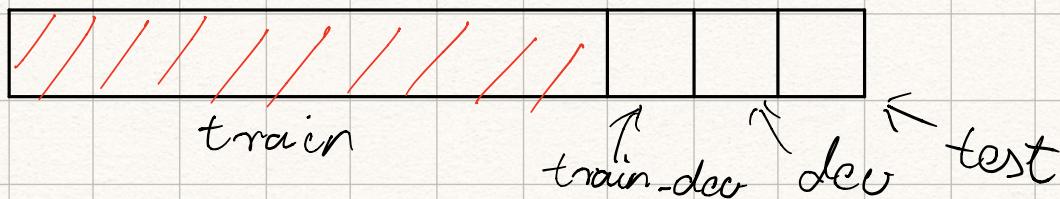


same distribution, but web will dominate



\Rightarrow will do better on real applications

Bias and variance with mismatched distribution



Human error $\sim 0\%$

Train error

Train-dev error

dev error

1%	↕
9%	
10%	

variance

1%	↕
1.5%	
10%	

data mismatch

10%	↕
11%	
12%	

avoidable bias

10%	↕
11%	
20%	

Bias + data mismatch