

## o Course overview

logistic regression

stochastic optimization

Data & parameter tuning

Deep networks

Regularization

Convolutional Networks

Embeddings

Recurrent models

## o supervised Classification

### o logistic classifier

linear classifier:  $WX + b = y \Rightarrow \begin{bmatrix} 3.0 \\ 1.0 \\ 0.1 \end{bmatrix} \xrightarrow{\sigma} \begin{bmatrix} 0.7 \\ 0.2 \\ 0.1 \end{bmatrix}$

$\begin{matrix} \uparrow & & \uparrow \\ \text{weights} & & \text{bias} \\ & \searrow \quad \swarrow & \\ & \text{trained} & \end{matrix}$

$$y \Rightarrow S(y_i) = \frac{e^{y_i}}{\sum_j e^{y_j}} \Rightarrow p = 0.7$$

$\begin{matrix} \uparrow & & \uparrow & & \uparrow \\ \text{logistic scores} & & \text{softmax func} & & \text{probabilities} \end{matrix}$

## One-hot Encoding

$$A = \begin{bmatrix} 0.7 \\ 0.2 \\ 0.1 \end{bmatrix} \begin{matrix} \rightarrow a \\ \rightarrow b \\ \rightarrow c \end{matrix} \Rightarrow a = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad b = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \quad c = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

## Cross-Entropy

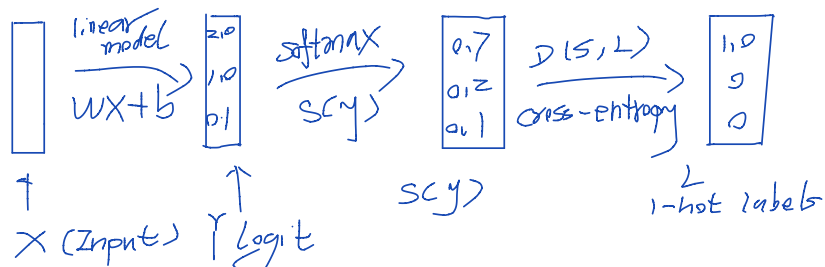
Compare classifier outcomes to one-hot encoding vectors to see how well the classifier does:

$$A = \begin{bmatrix} 0.7 \\ 0.2 \\ 0.1 \end{bmatrix} \leftrightarrow a = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \quad b = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \quad c = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

## Cross-entropy

$$\begin{matrix} \begin{bmatrix} 0.7 \\ 0.2 \\ 0.1 \end{bmatrix} & \xrightarrow{\quad} & D(S, L) = - \sum_i L_i \log(S_i) & \begin{bmatrix} 1.0 \\ 0 \\ 0 \end{bmatrix} \\ & \uparrow & \text{distance (entropy)} & \\ & \text{note: } D(S, L) \neq D(L, S) & & \end{matrix}$$

## Multinomial Logistic Classification



o Minimize Cross Entropy

$$\mathcal{D}(Scwx + b, L)$$

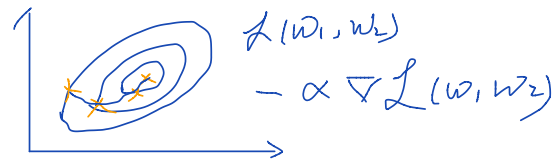
want: low entropy for correct class  
high entropy for incorrect class

Training Loss: average cross-entropy

$$\mathcal{L} = \frac{1}{N} \sum_i \mathcal{D}(Scwx_i + b, L_i)$$

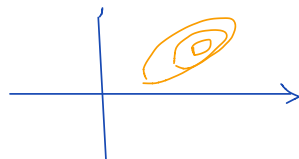
$\uparrow$  big sum       $\uparrow$  big matrix

o Gradient Descent

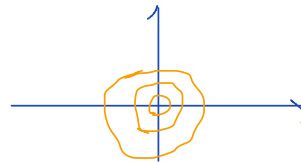


o Normalized inputs and initial weights

Badly conditioned



Well conditioned




zero mean, equal variance

normalization:  $\frac{x_i - \bar{x}}{\sigma}$

◦ Initialization of logistic classifier

$$\mathcal{L} = \frac{1}{N} \sum_i \mathcal{D}(S(Wx_i + b_i), L_i)$$

↑  
random gaussian  
distributed weights



◦ Optimization:

$$w \leftarrow w - \alpha \Delta_w \mathcal{L}$$

$$b \leftarrow b - \alpha \Delta_b \mathcal{L}$$

◦ Measuring Performance

Training  $\xleftrightarrow{\text{tune}}$  Test  $\longrightarrow$  unseen test Large Error!!!

You are giving your classifier the test set information through your eyes !!!

Solution: validation set

◦ Overfitting  $\longrightarrow$  Data size

The larger the data size, the less likely you will overfit

## o Stochastic Gradient Descent (SGD)

logistic regression: very hard to scale

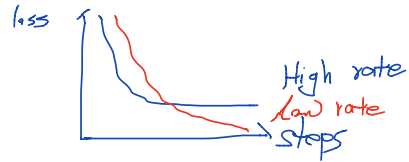
$L = \sum_{i=1}^N J_i$  All DATA! for gradient descent

## o Helping SGD

Running average: memory of history  
learning rate decay: convergence

## o Parameter hyperspace

learning rate tuning:



SGD "Black magic"

— Hyperparameters

- Initial learning rate
- learning rate decay
- momentum
- Batch size
- weight initialization

- Remember

Keep Calm And  
Lower Your  
Learning Rate!

- ADGRAD

only need to tune

- Batch size
- Initial weight.