

## Introduction:

Bias vs variance : a perfect decision  $y = f(x)$ ;  $\hat{y}$  train label, Bias - dis between average model and theoretical best

Simple model : high Bias, low variance Hyper-para : # of samples, feature dim d. Cross-validation : partition data into  $N$  subset, every time train on  $N-1$ ,

Complex model : low Bias, high variance. regularization :  $\lambda$  some parameter.

make use of small data set, but expensive.

variance - variability with different training samples.

Cross-validation : validate on  $N$ , rotate all  $N$  options, choose best configuration, remain

$$\begin{aligned} \text{MSE} &= E[(\hat{y}_0 - f(x_0))^2] = E[(\hat{y}_0 - E[\hat{y}_0] + E[\hat{y}_0] - f(x_0))^2] \\ &= E[(\hat{y}_0 - E[\hat{y}_0])^2] + E[(E[\hat{y}_0] - f(x_0))^2] + 2E[(\hat{y}_0 - E[\hat{y}_0])(E[\hat{y}_0] - f(x_0))] \\ &= E[(\hat{y}_0 - E[\hat{y}_0])^2] + (E[\hat{y}_0] - f(x_0))^2 \quad \text{for } E[\hat{y}_0 - E[\hat{y}_0]] = 0. \\ &= \text{Var}(\hat{y}_0) + \text{Bias}^2(\hat{y}_0) \end{aligned}$$

## Regression :

linear  $f(x; w) = w^T x$ , extend to non-linear  $f(x; w) = w^T h(x)$ ,  $h(x)$  non-linear.

OLS : total error  $= \sum_i (y_i - \hat{y}_i)^2 = \sum_i (y_i - \hat{y})^2$ ,  $\hat{y}$  prediction.

Matrix notation:  $(Hw - Y)^T(Hw - Y)$ ,  $H$ :  $N \times k$ ,  $w$ :  $k \times 1$ ,  $Y$ :  $N \times 1$

close form solution:  $L = (Hw - Y)^T(Hw - Y)$ ,  $L$  scalar, target  $\frac{\partial L}{\partial w} = 0$

$$\text{tr}(L) = L \Rightarrow \frac{\partial L}{\partial w} = \frac{\partial L}{\partial w} + (Hw - Y)^T(Hw - Y) = \frac{\partial L}{\partial w} + (w^T H^T H w - w^T H^T Y - Y^T H w + Y^T Y)$$

property:  $\text{tr}(AB) = \text{tr}(BA)$ .

$$\text{tr}(ABC) = \text{tr}(CAB) \dots \text{tr}(ABC)$$

$$d\text{tr}(wB) = d\text{tr}(Bw) = B$$

$$d\text{tr}(A^T B) = B$$

$$d\text{tr}(A) = d\text{tr}(A^T)$$

$$d\text{tr}(AxB^T) = AxB + A^TxB^T$$

$$d\text{tr}(A^TxB^T) = AB$$

$$\frac{\partial L}{\partial w} + (w^T H^T H w) = (H^T H w) + (w^T H^T H^T)$$

$$\frac{\partial L}{\partial w} + (w^T H^T Y) = \frac{\partial L}{\partial w} + (Y^T H w) = (Y^T H^T)$$

$$\frac{\partial L}{\partial w} + (Y^T Y) = 0$$

$$\frac{\partial L}{\partial w} = H^T H w + H^T Y = 0$$

$$H^T H w = H^T Y \Rightarrow w = (H^T H)^{-1} H^T Y.$$

Regularization : large parameter cause overfit.

Ridge:  $\hat{w}_{\text{ridge}} = \arg \min_w \frac{1}{2} \|w\|^2 + \frac{1}{2} \|y - (Hw - Y)\|^2 = \arg \min_w (Hw - Y)^T (Hw - Y) + \lambda \|w\|^2$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial w} + \text{tr}(w^T H^T H w - w^T H^T Y - Y^T H w + Y^T Y + \lambda w^T I w)$$

$$= 2H^T H w - 2H^T Y + 2\lambda I w = 0.$$

$$(H^T H + \lambda I) w = H^T Y$$

$$w = (H^T H + \lambda I)^{-1} H^T Y.$$

LASSO:  $\hat{w}_{\text{lasso}} = \dots + \sum_{k=1}^K w_k$

push more weight to 0.

allow feature selection.

no differentiable and close form.

## CNN:

$$\text{sigmoid: } f_{\text{sig}} = \frac{1}{1+e^{-x}}$$

$\frac{df_{\text{sig}}}{dx} = f_{\text{sig}}(1-f_{\text{sig}})$  (MLP)

$$\tanh: \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

(CNN)

diff MLP and CNN:

use convolution (height sharing), pooling (subsampling), no feature extraction as pre-processing, use tanh.

Receptive Field: region of input space affects a particular unit of the network.

Dilated conv: enlarge the receptive field. dilation = 2,  $\text{Tr}(I^2)$  input

[xi] conv: cross channel information aggregation, feature projection, dim reduction.

Depthwise conv: each channel has its own kernel, parameter reduced, but no cross channel information (combine xi conv)

group conv: split channels into group, still need xi conv. (ResNeXt)

DPM: sliding window, binary classification, HOG as feature, Linear SVMs.

2-CNN: extract region (selective search), get feature (CNN), classify (SVM). Mask R-CNN: Mask prediction, ROI align.

cls: 该框置量计算, 通过框的尺寸, 高度, 宽度, 多维特征值, 缩放, 增加置信度

Fast D-CNN: one-stage train, no sum, no warp/resize, fast RPN.

YOLO: use stride to decrease size, no pooling and full-connect, fast, simple algorithm, larger scale, bad at small object.

SSD: multi box and multi scale.

NMS: algorithm to choose the best bounding box.

## SVM:

optimal margin classifier

$$: h_0(x) = g(w^T x), w^T x \geq 0, \text{Predict } "1" \Rightarrow \text{if } y^{(i)} = 1, \text{ we hope } w^T x \gg 0$$

geometric margin (linearly separable)

$$\text{label } y \in \{0, 1\}, h_0(x) = \frac{1}{2} \text{ if } x=0, \text{ otherwise, } h_0(w) = g(w^T x + b)$$

no constraint  $x=0$

functional margin by  $(w, b)$   $\hat{w}$ : st. train sample is:  $\hat{w}^{(i)} = y^{(i)}(w^T x^{(i)} + b)$ , hope  $\hat{w}^{(i)} > 0$

normalize:  $(w, b) \Rightarrow (\frac{w}{\|w\|}, \frac{b}{\|w\|})$  to avoid cheating.

geometric margin of hyperplane ( $w, b$ ) wrt  $(x^{(i)}, y^{(i)})$  is  $r^{(i)} = \frac{(w^T x^{(i)}, y^{(i)})}{\|w\|}$  distance.

functional margin

optimal margin classifier choose  $w, b$  to max  $r$ .

max r. st.  $y^{(i)}(w^T x^{(i)} + b) \geq r, i=1, 2, \dots, n \Rightarrow \min_i \|w\|^2 \text{ s.t. } y^{(i)}(w^T x^{(i)} + b) \geq r$

$\|w\|^2 = \frac{1}{2} w^T w$

if choose  $\|w\| = \frac{1}{r}$ ,  $\Rightarrow \max_i y^{(i)}(w^T x^{(i)} + b) \geq 1$ .

suppose  $w = \frac{1}{r} \text{ diag}(w^T x^{(i)})$ ,  $\Rightarrow \min_i \frac{1}{2} \|w\|^2 \text{ s.t. } y^{(i)}(w^T x^{(i)} + b) \geq 1$ ,

$\Rightarrow \text{as min} \frac{1}{2} \sum_i \text{diag}(y^{(i)}(w^T x^{(i)} + b))$ ,  $y^{(i)}(\frac{1}{r} \text{ diag}(w^T x^{(i)})^T x^{(i)} + b) \geq 1$

equivalent to  $\max_i \frac{1}{r} \text{ diag}(y^{(i)}(w^T x^{(i)} + b)) \geq 1$ , s.e.  $\text{diag}(w^T x^{(i)} + b) \geq r$ . Dual optimization.

solve  $\alpha, b$ , compute  $h_0(w, b) = g(w^T x + b) = g(\frac{1}{r} \text{ diag}(w^T x) + b) = g(\frac{1}{r} \text{ diag}(x^T w) + b)$ .

kernel trick: write algorithm in terms of  $x^T x^T x^T$ , mapping from  $x \rightarrow \phi(x)$ , high dim., find  $k(x, z) = \phi(x)^T \phi(z)$ , replace  $x, z$  with  $k(x, z)$

$J(\theta) = C \sum_i [cost(\theta|x_i|) + (1-y^{(i)})[cost(\theta|x_i|)] + \frac{1}{2} \sum_j \theta_j^2]$

Support vector  $y^{(i)}(w^T x^{(i)} - b) = 1$  s.t.  $y^{(i)}(w^T x^{(i)} - b) \geq 1 - \varepsilon_i, \varepsilon_i \geq 0$ , or slack variable

for non-linearly:  $\min_{w, b} \frac{1}{2} \|w\|^2 + C \sum_i \varepsilon_i$  s.t.  $y^{(i)}(w^T x^{(i)} - b) \geq 1 - \varepsilon_i, \varepsilon_i \geq 0$ , or slack variable

目标: max margin, min mis-classified.

Lagrange duality.

**NLP:** Vector Semantics: words meaning characterized by words frequently appear close-by. distributional semantics.

Bag of words: word counts, document-term matrix. / word-word matrix: window of K nearby words  $P_{\text{PMI}}(l_1, l_2) = \max(\log_2 \frac{P(l_1, l_2)}{P(l_1)P(l_2)}, 0)$

word embedding: prediction based, cbow/skip-gram.

CBOW:  $w_u$ ,  $w_c$  center word,  $w_{bc}$  bc context word  $\text{CBOW} = \frac{\exp(\frac{1}{n} \sum_i \delta_i)}{\sum_i \exp(\frac{1}{n} \delta_i)}$  center C, center O ;  $P(w_c | w_{bc}, \dots, w_{bc}) = \frac{\exp(\frac{1}{n} \sum_i \delta_i)}{\sum_i \exp(\frac{1}{n} \delta_i) (V_{bc} + V_{oc})}$ ,  $\log P(w_c | w_{bc}) = \delta_i^T w_c - \log \sum_i \exp(\delta_i^T w_c)$

$\frac{\partial \log P(w_c | w_{bc})}{\partial w_c} = \frac{1}{n} \delta_i^T (w_c - \frac{1}{n} \sum_i \delta_i w_c) / (V_{bc} + V_{oc})$

Fast training, efficient usage of statistics / primarily used to capture word similarity, disproportionate importance given to large counts) hard to see

scales with corpus size, inefficient usage of statistics / generate improved performance on other tasks. Capture complex patterns beyond word similarity.

word vec: N维； input one-hot. linear hidden, softmax output.

skip-gram: give specific word, pre-ready word. CBOW: glo context, pre-target word. CBOW: learn to pre word by context, skip-gram: pre context by center word.

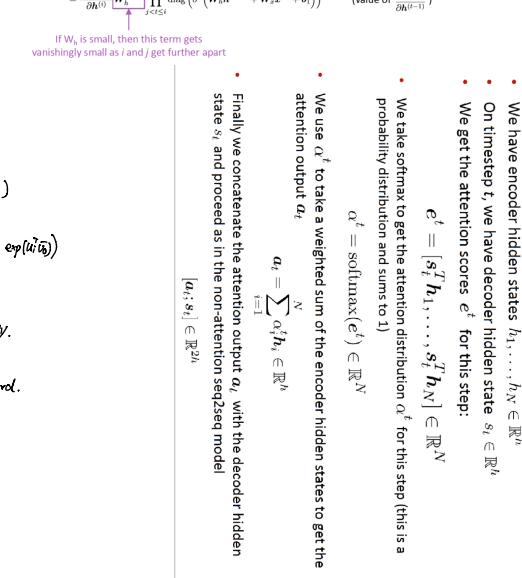
CBOW faster, skip-gram work well with small training data / huge vocab size, operation in softmax layer.

glove: word of co-occurrence probabilities can encode meaning components, only ratio, not counts.  $w_i \cdot w_j = \log P(i, j)$

Hierarchical softmax: tree, Log(LV) parameter  $|V|-1$  nodes.  $P(w_i | w_{bc}) = \sigma(u_i^T (w_{bc})) v_c \cdot \sigma(u_{bc}^T (w_{bc})) \cdot \sigma(u_{bc}^T (w_{bc})) v_c$

negative sampling:  $\nabla \theta = -h^T v_w + \log \sum_{i \neq w} P(w_i) \exp(h^T v_{w_i})$   $\nabla \theta \cdot \theta = \nabla \theta \cdot \theta = \sum_{i \neq w} P(w_i) \exp(h^T v_{w_i})$

task: language modeling, predict what words come next.



Vanish/Exploding Gradient

- Recall:  $h^{(i)} = \sigma(W_h h^{(i-1)} + W_x x^{(i)} + b_i)$
- Therefore:  $\frac{\partial h^{(i)}}{\partial h^{(i-1)}} = \text{diag}(\sigma'(W_h h^{(i-1)} + W_x x^{(i)} + b_i)) W_h$  (chain rule)
- Consider the gradient of the loss  $J^{(i)}(\theta)$  on step i, with respect to the hidden state  $h^{(j)}$  on some previous step j.

$$\begin{aligned} \frac{\partial J^{(i)}(\theta)}{\partial h^{(j)}} &= \frac{\partial J^{(i)}(\theta)}{\partial h^{(i)}} \frac{\partial h^{(i)}}{\partial h^{(i-1)}} \dots \frac{\partial h^{(j)}}{\partial h^{(j-1)}} \quad (\text{chain rule}) \\ &= \frac{\partial J^{(i)}(\theta)}{\partial h^{(i)}} \frac{\partial h^{(i)}}{\partial h^{(i-1)}} \dots \frac{\partial h^{(j)}}{\partial h^{(j-1)}} \text{diag}(\sigma'(W_h h^{(i-1)} + W_x x^{(i)} + b_i)) \quad (\text{value of } \frac{\partial h^{(i)}}{\partial h^{(i-1)}}) \end{aligned}$$

If  $W_h$  is small, then this term gets vanishingly small as  $i$  and  $j$  get further apart

- We use  $\alpha_t^i$  to take a weighted sum of the encoder hidden states to get the attention output  $a_t$  for this step: We take softmax to get the attention distribution  $\alpha_t^i$  and sums to 1. Finally we concatenate the attention output  $a_t$  with the decoder hidden state  $s_t$  and proceed as in the non-attention seq2seq model:  $[a_t; s_t] \in \mathbb{R}^{2L}$ .
- $\alpha_t^i = \text{softmax}(e^i) \in \mathbb{R}^N$

Introduction:

**Active learning:** the ability to ask queries about the world based upon the past queries and response

**Step:** Start from a pool of unlabeled data ( $W$ ) ; pick a few points randomly and get labels by oracle ; **Repeat :** Fit a classifier to labels seen so far pick best to see labels

**Stream-based Selective Sampling**: Using when obtain unlabeled instance is cheap; query can first be sampled from actual distribution, learned decide whether request labeled best; (close to boundary) ; most uncertain? most likely to  $\downarrow$  Overall uncertainty).

**Pool-based Sampling:** Assume small set of labeled data and large pool of unlabeled  $\mathcal{U}$ , former supplied with set of unlabeled examples to select pretrain.   
 *pool-based sampling* / sequential active learning.

**Query Strategy Frameworks:**  
**Query-Driven Sampling**: desire small set of labeled data and large pool of unlabeled data, learner supplied with set of unlabeled examples to select from  
**Uncertainty Sampling**: queries the instances least certain how to label (close to boundary), 3 strategies: least confident, margin sampling, entropy.

**Query-By-Committee:** Given a committee of models trained on labeled set  $L$ , but competing each other, each works for label of query candidates. Pick the most disagreement one, goal is min version space.

$\hat{X}^u = \arg \min_{\hat{X}} \frac{1}{n} \sum_{i=1}^n \log \frac{p_{\text{obs}}(x_i)}{p_{\text{pred}}(x_i)}$ , where  $n$  is number of vertices,  $x_i$  is community membership vector.

Empirical gradient w.r.t.  $\theta$  (Eqn. 1) is called gradient-based training.

$\times \hat{y}_n = \arg \min_{\beta} \sum_{i=1}^n P(y_i | x_i) = \left[ \sum_{i=1}^n \beta^* + \beta y_i - \log(1 + e^{(\beta^* + \beta y_i)}) \right]$ ,  $\beta \rightarrow \infty$  as new model often  $\rightarrow$  fits added and retrain.

$X^*$  is defined as model's "quasi-the most informative and representative" (invariant regions of input space)  $\arg \max_{\mathbf{x}} f(\mathbf{x}) \times \left( \frac{1}{k} \sum_{i=1}^k \sin(x_i, x_{i+1}) \right)^{\beta}$ . False information of  $\mathbf{x}$  based on some base query strategy A,  $\beta$  controls relative importance.

its change  
iance .  
of  $\Sigma$  .

den no p  
as co-v  
and valis

variance ).  
 $i$ . Step

method (minimise each layer's error)

is big error  
 $\Sigma \delta_i u_j$   
 step.  
 to Class  
 purity  
 $\frac{K}{n} = \frac{1}{2} \sum_{j=1}^n \delta_j$   
 error  
 $\frac{1}{n} \sum_{j=1}^n \chi_j X$   
 $= \frac{1}{1+e^{-x}}$   
 $\times j^{th}$  new  
 neuron in  
 $\frac{\partial \delta_j}{\partial x} = \frac{\partial \delta_j}{\partial \chi_j} \frac{\partial \chi_j}{\partial x}$   
 classif  
 some of  
 to get  
 feature  
 as acc  
 converg  
 due to re

of all pairs  
class closer  
 $\frac{x_i}{m}$ ,  $\frac{x_j}{m}$   
 $\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^m$   
 $(x_i - \bar{x})^T (x_j - \bar{x})$   
below class  
S.  
sigmoid function  
weight function  
 $f(x) = \frac{1}{1 + e^{-w^T x}}$   
 $w = (w_0, w_1, \dots, w_n)$   
 $b = b_0$   
 $h_\theta(x) = \theta^T x$   
loss function  
steeper gradient  
gradient descent  
sparse matrix  
fastest  
(book), Pearson

er, but exp  
small num  
m.  
(1). Average  
 $\bar{x} = \frac{1}{n} \sum x_i$   
clustered  
water.  
 $j = \frac{1}{2} \cdot \frac{1}{2}$   
ability of  
application  
normal basis  
 $x_1^2, x_2^2$   
 $X_1, X_2$   
eigenvalue  
 $\frac{\sigma_x^2 - \sigma_y^2}{\sigma_x^2 + \sigma_y^2}$   
er., W.  
in  
ex. feature  
 $\frac{\partial}{\partial w_{j+1}} \dots \frac{\partial}{\partial w_k} f'(K+1)$   
 $\frac{\partial}{\partial b}$   
loss entropy  
gross margin  
subsequence  
ter,  
g rate  
rep. learn  
Scaling  
squared

different lines  
decisions for  
a develop-  
ment problem.  
Assign  
ation  
/ divisive  
overlap  
 $\Sigma_{j=1}^n$   $(\bar{x}_j; \bar{y}_j)$   
is, prob-  
a some  
1, ortho-  
 $\bar{x} = \frac{1}{n} \bar{\mathbf{x}}$   
in row of  
highest  
int f(z);  
at 2nd  
as activa-  
in rule,  
 $f(w, b)$ ,  
(b)  
ji  $\Sigma$   
 $(W(K)^T$   
 $|z|_2, \dots$   
use  
station :  $\mathbb{P}$   
age over  
h param  
, learn  
problem  
s, better  
ents are

group of samples ( $k$ )  
 (clusterings)  
 single-link  
 cluster C<sub>i,j</sub>  
 $(x_i, c_i^2)$   
 multiple clusters with  
 $\hat{x}_i = \frac{1}{k} \sum_{j=1}^k x_{ij}$

noise (iter)  
 ly divide  
 family of  
 same pair  
 as cluster  
 $\{x_i\} = \{x_j\}$  dis-  
 tance not de-  
 creasing  
 ..., allow  
 $x_i \in \{x_j\}$   
 ty of c  
 to lat  
 fections  
 basis  
 $i = 1, \dots$   
 $n_{\text{ex}}$ .  
 ents:  
**work**  
 $\Omega = \{f^{(k)}\}_{k=1}^K$   
 soft plus  
 $+ W_k x_k$   
 $b^{(k)}$ ,  $a^{(k)}$   
 $b^{(k)} - a^{(k)} =$   
 $s_i, s_i^{(k)}$   
 $\Rightarrow f = J(\Omega)$   
 $J(\Omega)$   
 $f^{(k)}$   
 $\delta_i$   
 $k, T$ ,  
 hidden  
 avoid 0  
 , ill-  
 um: E  
 of d  
 scaling  
 ent, br  
 us grav  
 enge che  
 I mousi

+ random  
 i, f<sub>i</sub>, f<sub>i+1</sub>, f<sub>i+2</sub>, f<sub>i+3</sub>, f<sub>i+4</sub>  
 location = i  
 $\hat{A}_{ij}^{\text{true}} = \sum_{k=1}^{n-1} P_j$   
 the Probabilistic  
 : Computer  
 true P  
 projection:  
 $\dots, x_i, x_{i+1}, \dots$   
 $X \in$   
 al Comp  
**I Need**  
 :  $h_{\alpha\beta}$   
 $\max(0, z) =$   
 $x_i + h_{i+1}^{(1)}$   
 $h_{i+1}^{(1)} =$   
 $W^{(1)} = W^{(1)}$   
 $h_{i+1}^{(1)} = 0$   
 $d_{ij} = d_{ij}^{(1)}$   
 $|h_{\alpha\beta}| =$   
 $(A_1^{(1)})^{-1}$   
 $= d_{ij}^{(1)}$   
 $W^{(1)}$   
 number of  
 sian to  
 checking  
 Momen  
 direction  
 tomatrical  
 & adjust  
 all pre  
 moving  
 Experiment

- back propagation  
 $\Rightarrow \frac{\partial J}{\partial W_{ij}}$   
 $W_{ij}^{(l)}$   
 $\delta^{(l+1)} = \text{choose } \tau$   
 $\text{regularization}$   
 $\text{gradient}$   
 SGD with  
 Stab.  
 Ada:   
 Smo.  
 over.  
 RMS Prop  
 Adam

hierarchical	
2. Distance	
Agglomerative	
K-means	
hard assign	
Gmm	
Evalu	
PCA	
Li	
X	
alg	
M	
a pe	
relu	
Gr	
Form	
Back	
W	
δb	
1. δ	
2. δ	
2. δ	
char	
reg	
Crat	
Sc	
Ad	
RMS	
Adv	

$$f_{\text{conv}}(\cdot) = \rho_{\text{conv}}(\cdot)$$

$$I(X;Y) = H(X) - H(X|Y)$$

$$H(X) = \sum_x P_x(x) \log P_x(x)$$

$$H(X|Y) = \sum_y P_y(y) \left( \sum_x P_{xy}(x|y) \log P_{xy}(x|y) \right).$$

$$H(X;Y) = H(X) - H(X|Y)$$

$$H(X) = -\sum_x P_X(x) \log P_X(x)$$

$$H(X|Y) = \sum_y P_Y(y) \left( -\sum_{x|y} P_{X|Y}(x|y) \log P_{X|Y}(x|y) \right).$$

<p>• Proposed by Cho et al. in 2014 as a simpler alternative to the LSTM.</p> <p>• On each timestep we have input <math>x^{(t)}</math> and hidden state <math>h^{(t)}</math> (no cell state).</p>
<p>hidden state: controls what part of previous hidden state are preserved</p>
<p><b>Reset gate</b> controls what parts of previous hidden state are used to compute new content</p>
<p><b>New hidden state element</b>: need gate selects useful parts from past state, this and current input to compute new hidden content</p>
<p><b>Hidden state: update gate</b> simultaneously controls what is kept from previous hidden state, and what is updated to new hidden state content</p>
<p>How does this solve vanishing gradient? Long LSTM, GRU makes it easier to get into long-term (e.g. by setting update gate to 0)</p>