

**Supervised (Inductive) Learning:** Train + labels (desired output) e.g. classification, regression, obj detection, semantic segmentation

**Unsupervised:** Train (no desired output) e.g. clust, dim reduction

**Semi-supervised:** Train + few desired output

**Self-supervised:** Pretext task (e.g. predict if rotated image), learn the latent features (output labels) intrinsically from data

**RL:** Rewards from seq actions

**Evaluation:** Acc (must be class specific in case of imbalanced data)

E.g. If only 30% 1 & predict all 0, we have 70% acc! totally biased

Predicted

Actual			Recall/ Sensitivity: $\frac{TP}{TP+FN}$	Specificity: $\frac{TN}{TN+FP}$
	+ve	-ve		
	+ve	TP		
-ve	FP (Type I Error)	TN		
Precision:	$\frac{TP}{TP+FP}$	-ve Pred Val: $\frac{TN}{TN+FN}$	Acc:	$\frac{(TP+TN)}{(All)}$
F1:	(Precision * Recall)/(Precision + Recall)			

**MSE (L2 Loss) for regression:**

$$MSE = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

**Receiver Operating Curve (ROC):**  
TPR vs FPR (plot against all threshold, AUC = 1 best)

**Information:** (Generally) Measured in bits ( $\log_2$ ): Amt of "surprise" arising from a given (stochastic) event E:  $I(E) = -\log(P(E))$

E.g. Three boys: Tom (20%), Bob (1%), Sam (79%)

Bob =  $-\log_2(0.01) = 6.6$ , high surprise vs

Sam =  $-\log_2(0.79) = 0.34$ , low surprise

**Entropy:** Avg lvl of information inherent in possible outcomes:

Higher Entropy signifies more uncertainty in what action an agent would choose:  $H(X) = -\sum_i P_X(x_i) \log_b P_X(x_i)$

$$Entropy = -(0.20 \log_2(0.20) + \dots + 0.79 \log_2(0.79)) = 0.26$$

**KL-Divergence (relative entropy):**

Expression of Surprise, diff btw 2 prob distribution P & Q:

Under the assumption P & Q are not close, KL divergence high (surprise). If close then KL divergence low

**Likelihood ratio (n indep samples):**  $LR = \prod_{i=0}^n \frac{p(x_i)}{q(x_i)}$

(how much more likely samples came from P dist than Q dist)

**Log likelihood ratio:  $LLR = \sum_{i=0}^n \log\left(\frac{p(x_i)}{q(x_i)}\right)$**

**Expected LLR:**

(Avg, how much each sample gives evidence of P(x) over Q(x))

$$D_{KL}(P \parallel Q) = \sum_{i=0}^n p(x_i) \log\left(\frac{p(x_i)}{q(x_i)}\right)$$

$$D_{KL}(P \parallel Q) = \sum_{i=0}^n p(x_i) \log(p(x_i)) - \sum_{i=0}^n p(x_i) \log(q(x_i))$$

Entropy (Info content of P) - IC of Q weighted by P

If P is the "true" distribution, then the KL divergence is the amount of information "lost" when expressing it via Q.

**Cross-Entropy (min unlike max likelihood):**

Cross entropy is, at its core, a way of measuring the "distance" between two probability distributions P and Q. (Entropy on its own is just a measure of a single probability distribution)

CE = Combination of the entropy of the "true" distribution P and the KL divergence between P and Q:

$$H(p, q) = H(p) + D_{KL}(p \parallel q) = -\sum_{i=0}^n p(x_i) \log(q(x_i))$$

E.g. Predict one-hot vector: True {0, 0, 1, 0, 0}

Output (of softmax layer): Predicted {0.01, 0.02, 0.75, 0.05}

Predicted class (argmax): 2 (zero index)

$$-(0 \cdot \log(0.01) + 0 \cdot \log(0.02) + 1 \cdot \log(0.75) + 0 \cdot \log(0.05))$$

**Cross-entropy:**

Avg # of total bits needed to encode data from a source with dist p when we use model q

**KL-Divergence:**

Avg # of extra bits needed to encode the data, when using dist q instead of the true dist p.

**Logistic Regression:**

Given an input space X, the goal is to predict for every sample  $x \in X$  to which class it belongs. For 2 class classification, the g predict output space  $Y = \{0, 1\}$  or  $\{-1, +1\}$

Linear:  $f(x) = w \cdot x + b$ , unbounded space

Need to map  $(-\infty, \infty)$  to  $(0, 1)$  st 0 map to 0.5

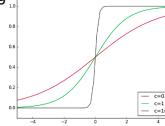
Use logistic sigmoid fn  $s(a)$ :

$$s(a) = \frac{\exp(a)}{1 + \exp(a)} = \frac{1}{\exp(-a) + 1} \exp(a) = \frac{1}{\exp(-a) + 1}$$

Note if multi class (k) class, use Softmax:

$$\sigma(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \text{ for } i = 1, \dots, K \text{ and } \mathbf{z} = (z_1, \dots, z_K) \in \mathbb{R}^K$$

$$h(x) = s(f_{w,b}(x)) \in (0, 1) \text{ aka prob(sample x class label } y = +1)$$



**Cross-Entropy loss (bin class):**

**CE Loss:** (1 sample:  $(x_i, y_i)$ )

$-\bar{y}_i \log(p_1(x_i)) - (1 - \bar{y}_i) \log(1 - p_1(x_i))$  where  $p_1(x_i) = 1 - p_2(x_i)$

**Avg CE Loss:** (dataset: n indep samples)

$$\frac{1}{n} \sum_i -\bar{y}_i \log(p_1(x_i)) - (1 - \bar{y}_i) \log(1 - p_1(x_i))$$

Opti problem (convex) to be solved is (finding best parameter w min CE loss):

$$(w^*, b^*) = \operatorname{argmin}_{w,b} (\text{avg CE Loss})$$

$$\text{where } p_1(x_i) = h(x_i) = \frac{1}{\exp(-w \cdot x_i - b + 1)} = s(w \cdot x_i)$$

$$\nabla_w L = \nabla_w \left( (-1) \cdot \sum_{i=1}^n y_i \log(s(w \cdot x_i)) + (1 - y_i) \log(1 - s(w \cdot x_i)) \right) = \sum_{i=1}^n x_i (s(w \cdot x_i) - y_i) = \sum_{i=1}^n x_i (h(x_i) - y_i)$$

**Data Set:** (No best split) **Training:** run your learning algo

**Development:** tune parameters **Test:** evaluate perf

**Common Error:**

1) Data mismatch (e.g. diff. Resolution img)

2) Bias and variance

**Bias:** error rate on the training set

**Variance:** how much worse test set error than train (e.g. 15% error (85% acc) on train, bias = 15%. But 16% error on test, variance = 1%)

**Set Distribution:**

If all same dist, then optimising for web imgs mostly here, which is diff from target dist.

100k imgs (from Web) 8k imgs (from Target dist)

Train (104k imgs) Dev (2k imgs) Test (2k imgs)

Instead make test & dev from the same dist such that optimizing to perform well on the target dist.

100k imgs (from Web) 4k imgs 4k imgs

Train (104k imgs) Dev (2k imgs) Test (2k imgs)

However, the training dist is now diff from the dev/test distribution: it will take longer & more effort to optimize the model. More importantly, cannot easily tell if the classifier error on the dev set relative to the error on the train set is a variance error, a data mismatch error, or a combination of both.

You can use a **bridge set (training data not used for training)** to trouble shoot where the error comes from.

100k imgs (from Web) 4k imgs 4k imgs

Train (102k) Bridge (2k) Dev (2k) Test (2k)

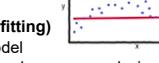
Error 2% 3% 10%

Error 8% (btw dev & train): 1% var + 7% data mismatch Hence our classifier performed well on data from the same distribution but poorly when from a diff distribution, (data mismatch problem). Else if higher var & lower data mismatch error, implies overfitted to train set.

**Rules of thumb:**

**High avoidable bias: (underfitting)**

- Increase the size of your model
- Modify the input features based on error analysis
- Reduce or eliminate regularisation (reduces bias but increases variance)
- Modify model architecture



**High variance: (overfitting)**

- Add data to your training set (e.g. by data augmentation)
- Add regularisation
- Add early stopping
- Feature selection to decrease # of features (may increase bias)
- Decrease the model size (May increase bias)
- Modify input features based on error analysis



**Prediction Error == Bias<sup>2</sup> + Variance + Irreducible Error**

**Total Error** = **Variance** + **Bias<sup>2</sup>**

**Low bias, low variance**

**Model Complexity**

**Optimum Model Complexity**

**Desired perf**

**Dev error**

**Training error**

**Desired**

**Dev Error**

**Training error**

**Desired**

**Optimum Model Complexity**

**Desired perf**

**Dev error**

**Training error**

**Desired**

**Dev Error**

**Training error**

$\hat{y} = g(w_0 + x^T w)$  linear decis<sup>1</sup> no matter network size  
 Nonlinear  $g(a) \Rightarrow$  approx of arbit complex f<sup>2</sup>

- i) Threshold:  $[a > 0]$
- ii) Sigmoid  $\sigma = \frac{1}{1+e^{-a}}$
- iii) ReLU:  $\max(0, a)$

XOR: simplest linearly inseparable problem

a) Reg.: Linear  $(-\infty, +\infty)$  / ReLU  $(0, +\infty)$   $\nabla$  MSE loss

b) Bin class.:  $\sigma(0, 1)$   $\nabla$  BCE loss

c) Categorical: Softmax  $(0, 1)$   $\forall$  class,  $\sum_i p_i = 1$   $\nabla$  CE loss

d) Multiclass cat:  $\rightarrow$   $\nabla$  BCE loss  $\nabla$  class

Empirical loss: Total loss over entire data

$$J(W) = \frac{1}{n} \sum_{i=1}^n \ell(f(x^{(i)}, w), y^{(i)}) \quad \text{MSE: } \frac{1}{n} \sum_{i=1}^n (y^{(i)} - f(x^{(i)}, w))^2$$

$$\text{BCE: } \frac{1}{n} \sum_{i=1}^n [y^{(i)} \log(f(x^{(i)}, w)) + (1-y^{(i)}) \log(1-f(x^{(i)}, w))]$$

Training: Find  $W^* = \underset{W}{\operatorname{argmin}} J(W)$ , lowest loss

i) Init weights rand  $\sim N(0, \sigma^2)$

ii) Loop until converge:  $\nabla W = \frac{\partial J(W)}{\partial W}$ ,  $W \leftarrow W - \eta \nabla W$

iii) Return  $W$ 

- Small  $\eta$ : converge quickly (local optima)  $\nabla$  Use adagrad!
- Large  $\eta$ : unstable & diverges

- Easy compute (implement<sup>1</sup> & understand<sup>2</sup>)

- May trap at local minimum (need convex f<sup>2</sup> too)

- Update weights only after calculating  $\nabla$  on entire dataset (too large  $\Rightarrow$  years to compute + large mem)

SGD: Update model param more freq.

- Converges in less time  $\nabla$  may get new minima's

- But high var in model param (may shoot up even after reaching global min)

- To get some convergence as GD, need slowly  $\eta$  to  $\eta$

Pick Batch B of data  $\nabla$  Fast compute & allow // compute better est. of true  $\nabla$  than 1 datapoint, smoother converge, allow larger  $\eta$

Adagrad: Changes  $\eta$  if param  $\nabla$  t  $\nabla$  RMS Prop

- No need tune  $\eta$  manually & trainable on sparse data

-  $\eta$  always  $\downarrow \rightarrow$  slow training

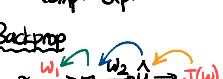
- Computationally exp (need 2nd order der)

ADAM: momentum of 1st+2nd order

- Very fast converge + rectifies decaying  $\nabla$  & high var

- Comp. exp.

$$\frac{\partial J(w)}{\partial w} = \frac{\partial J(w)}{\partial g} \cdot \frac{\partial g}{\partial z_1} \cdot \frac{\partial z_1}{\partial w}$$



Generative Modeling: Take input train samples from some dist and learn a model that represents that dist (Debiasing)

- Capable of uncovering underlying features in a dataset  
 - Detect outliers in the dist for training

AE: Encoder: Unsupervised learning of a lower dim feature space (under complete),  $z$  (salient features + easier to work with) representation from unlabelled train data

Decoder: Use  $z$  to reconstruct the obs  $\hat{x}$ . Note no labels used.  $\square z \square$

$L(x, \hat{x}) = \|x - \hat{x}\|^2$

- Form of compression, smaller  $z \rightarrow$  larger training bottleneck
- Bottleneck hidden layers forces network to learn compressed latent rep. while reconstruct loss forces  $z$  to capture as much information of the data

Regularizer AE: Overcomplete case with a loss function that encourages the model to have other properties beside copy-paste

- sparsity rep, smallness of rep derivatives (similar img, sim lv), robust to noise/missing inputs

Sparse AE:  $L(x, g(f(x)) + \Omega(z)$  decoder, encoder output + sparse penalty

Denoising AE: Add noise to input image, AE learns to remove noise

Bag word (Bow): Col: sentences  $\nabla$  Bin 0/1  $\rightarrow$  sparse Row: all words  $\nabla$  1-hot encode  $\uparrow$

Term-freq-Inv doc freq (TF-IDF): rare words carry more meaning

$$TF = \frac{\# \text{appearances for a term}}{\text{Total # words in doc}} \quad \left\{ \begin{array}{l} \text{TF-IDF} = \\ \text{TF} * \text{IDF} \end{array} \right. \\ \text{IDF} = \log(\# \text{docs} / \# \text{docs term appear in})$$

Loc-sparse rep: cooccurrence models (SVD)  
 vector embedding models (word2vec)

Word embeddings: each word = pt in n-dim space,  
 (1-layer) rep by vector of len n ( $\sim 100-300$ )

i) Cont-BOW: Prod center word from sum context  $\geq -$

ii) Skip-gram: Prod sum context words from given center word

- Max likelihood of context, given Raw word

e.g. I love pizza,  $P(I|pizza) \nabla P(\text{love}|pizza)$

Train word2vec: (weights of layer)  $\nabla$   $\times$  prepare order!

- Decide window size & embedding size & vocab size

- Training abg: Negative sampling (bin LR clas)

- Fast accurate model that capture semantic content:

"Gensim" lib, pretrained w.e.

- Interested in the 1 hidden(latent rep) layer, not the output

RNN: Output:  $\hat{y}_t = W^T h_t$  Input:  $x_t$  vector

Update hidden state:  $h_t = \tanh(W_h h_{t-1} + W_x x_t)$

i) Many values  $> 1$  (Exploding  $\nabla$ )  $\nabla$  clipping

ii) Many values  $< 1$ : Small # mult,  $\nabla$   $\uparrow$  smaller, bias param unable to capture short term dependencies

a)  $g(a)$ : ReLU, prevent f<sup>1</sup> from shrinking  $\nabla$  when  $a > 0$

b) Weight Init: Weights  $\rightarrow$  II matrix  $\nabla$  prevent weights Bias  $\rightarrow 0$  shrink to 0

c) Using gates to track & control info thru time

LSTM: - Forget irrelevant info of prev state

$$f_t = \sigma(W_f [h_{t-1}, x_t] + b_f)$$

- Store relevant new info into cell state

$$i_t = \sigma(W_i [h_{t-1}, x_t] + b_i) \quad \tilde{C}_t = \tanh(W_c [h_{t-1}, x_t] + b_c)$$

$$- \text{Update cell state selectively } C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

- Output gate: control what to send to next t

$$O_t = \sigma(W_o [h_{t-1}, x_t] + b_o) \quad h_t = O_t * \tanh(C_t)$$

GRU or attention networks (transformer)

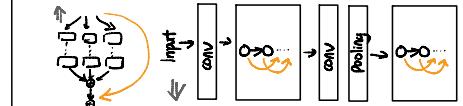
Transfer Learning (Opti): Higher start, slope & asymptote  
 - Identify related task with abundant data, reuse pretrained model

E.g. LeNet 5, AlexNet, VGG-16, Inception (GoogleLeNet)

- More layers -ve effect, degradation problem

ResNet: residual blocks at intermediate layers

ResNext: "Split-transform-merge" + residual block connection



DenseNet: Each layer's feature map concatenated to the input of every successive layer within dense block

Generative Adversarial Networks (GAN):

Generator: noise  $\rightarrow$  init of data to try trick discriminator

Discriminator: Tries identify data from fakes generated by G

noise  $\square Z \square G \square D \square$  2nn compete

Data: Real  $\rightarrow$  D use as +ve e.g., Fake  $\rightarrow$  D use as -ve e.g.

Alternative training:

D Training: D ignores G loss & uses its own loss

- Sample random noise to produce G output

- classify & calculate D loss

- BP. through D & G weight, only update G weights

G Training: Opposite

- G loss penalizes G for failing to fool D

If G perfect, D acc = 50%, D feedback gets less meaningful

overtime, may collapse, sin ce GAN not stable, but fleeting convergence

Minimax loss:  $\min_{\theta_g} \max_{\theta_d} [E_{X \sim P_{\text{data}}} \log D_{\theta_d}(c) + E_{Z \sim P(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z)))]$

D output real

D output fake

(CE loss btw real & fake dist)

- D max probability real & output prob fake low

- G no influence on 1st term, tries min 2nd term (1-...)  $\rightarrow$  get high output for fake to fool D

- Can get stuck at early stages of training,

.. use Modified minimax loss: max log  $D(G(z))$  instead of min  $(1 - D(G(z)))$

Deep Conv GANs: Uses leaky ReLu (-ve slope before  $x = 0$ ), sin ce ReLu is non-0 centered, non diff at 0 & dying problem

Conditional GANs: same feed label Y(paired data)

Used for object transfiguration

a) pix2pix(cGAN with original img as condition)

- Paired data, D compares  $G(x)$  and  $y$  given  $x$  aka generate img in different style from original img

- To make img less blurry use L1/L2

$$G* = \arg \min_{G} \max_{D} L_{cGAN}(G, D) + \lambda L_{L1}(G)$$

Cycle GAN (Unpaired data): Domain transformation

- Uses adversarial loss + inverse mapping (sin ce mapping highly constrained). Introduces a cycle consistency loss to push  $F(G(x)) \approx X$

- Not just img, e.g. waveform, but not spectrogram

- Used to convert accented speech