

Basics

(n choose p) = n! / (p!(n-p)!), (x+y)^n = sum\_{p=0}^n x^p y^{n-p}, x, y in R
(a+b)^3 = a^3 + 3a^2b + 3ab^2 + b^3

dim P\_p(R^d) = (d+p choose p) for all p in N\_0, d in N with lead. order
p -> inf: dim P\_p(R^d) = O(p^d); d -> inf: dim P\_p(R^d) = O(d^p)

f(x) = f(x\_0) + J\_f(x\_0)(x - x\_0) + o(||x - x\_0||)

(a c d) (c d) = (ac ad bc bd)

- grad\_x ||x||\_2^2 = 2x, grad\_x b^T x = grad\_x x^T b = b
grad\_x (b^T A x) = A^T b
grad\_x (x^T A x) = (A + A^T) x if A sym. = 2A x
grad\_x (c^T X b) = c^T b
grad\_x (c^T X^T b) = b c^T

LSE with Squared Matrix A^{n x n}

- Ax = 0, if det(A) != 0 => x = 0, else (n-r) free variables
Ax = b, if det(A) != 0 => unique solution
else { (n-r) free variables, Rank(A) = Rank(A|b) = r < n
no solutions Rank(A) != Rank(A|b)

LSE with General Matrix A^{m x n}

- Ax = b, if Rank(A) < Rank(A|b) => no solutions
else if Rank(A) = Rank(A|b) = r => { (n-r) free variables, r < n
Unique solutions r = n

Supervised Learning

- f: R^n -> R cont diff and x\* stationary point => grad f(x\*)^perp = 0
x is a local minimum => grad f|x = 0 and Hf|x >= 0
f: R^n -> R continuous is convex: for all t in (0, 1)
f(tx + (1-t)y) <= tf(x) + (1-t)f(y)
strictly convex if < holds, concave if > holds, Hessian > 0

Composition: X subset R^n convex, f\_1, ..., f\_s: X -> R convex

- sum\_{i=1}^s lambda\_i f\_i is convex for all lambda\_i in R\_{>=0}
2. max{f\_1, ..., f\_s} convex
3. f\_1(f\_2(x)) convex for f\_1 being non-decreasing

(Batch) GD: w\_{t+1} = w\_t - eta grad L(w\_t) with w^{(0)} picked arbitrary

- high momentum -> long-lasting U-turn, oscillation down
loss should be decreasing, convex in its argument & diff-able
only exponential and logistic suffice
L(w^{t+1}) = L(w^t + eta v) = L(w^t) + eta grad L(w^t) dot v + o(eta)
||w^{t+1} - w^t||\_2^2 = ||-eta grad L(w^t)||\_2^2 = eta^2 ||grad L(w^t)||\_2^2

SGD at each iteration, decrease of loss not guaranteed

Minibatch GD batch size large => loss plot smoothens out.

Average gradients from diff. images are expected to lead in the optimal direction in the weight space

Batch Normalization standardize only batches of activation units v\_i at t\_0 and training. Regularize by limit magnitude of v\_i,

avoids internal covariate shift, vanishing/exploding gradient -> larger learning rates, faster convg.

Dropout large weights often come as a result of overfitting and dropout helps by randomly dropping weights during training

Depth reduction multiplicative nature of the chain rule -> small network depth prevents vanishing/exploding gradients

Regression Objective: w\* := arg min\_w L(w)

y = w^T x + epsilon => y = X w\* + epsilon

Ordinary Least Squares with closed form solution

w\* = (X^T X)^-1 X^T y Theta(nd^2 + d^3)

L(w) := sum\_{i=1}^n (y\_i - w^T x\_i)^2 = ||X w - y||\_2^2

grad\_w L(w) = -2 sum\_{i=1}^n (y\_i - w^T x\_i) \* x\_i = 2 X^T (X w - y) ~ Theta(nd)

Ridge Regression with closed form solution. Lasso uses lambda ||w||\_1

w\* = (X^T X + lambda I)^-1 X^T y

L(w) := sum\_{i=1}^n (y\_i - w^T x\_i)^2 + lambda ||w||\_2^2

grad\_w L(w) = -2 sum\_{i=1}^n (y\_i - w^T x\_i) \* x\_i + 2 lambda w

Logistic Regression P[Y = y | x] = 1 / (1 + exp(-y w^T x))

L(w) := log[1 + exp(-y w^T x)]

grad\_w L(w) = 1 / (1 + exp(y w^T x)) (-y x)

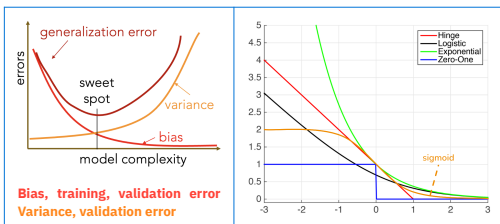
Regularized Logistic Regression

- L2: min\_w sum\_{i=1}^n log[1 + exp(-y w^T x\_i)] + lambda ||w||\_2^2
GD Step: w <- w(1 - 2 lambda eta) - eta grad\_w L(w)

Kernelized Logistic Regression

- alpha = arg min\_alpha sum\_{i=1}^n log[1 + exp(-y\_i alpha^T K\_i)] + lambda alpha^T K alpha
Classification P-hat(y|x, alpha) = 1 / (1 + exp(-y sum\_{j=1}^n alpha\_j k(x\_j, x)))

Cross-entropy loss l(W; x, y) = -log sum\_{y'} exp(f\_{y'}) / exp(f\_y)



Classification Objective: w\* := arg min\_w l(w; D)

0/1 Loss: l\_{0/1}(w; y\_i, x\_i) = { 1 if y\_i != sign(w^T x\_i), 0 else

Perceptron: algorithm uses together with SGD

- lin. separable <=> exists lin. separator (not necessarily optimal)

l\_p(w; y\_i, x\_i) := max(0, -y\_i w^T x\_i)

grad l\_p(w; y\_i, x\_i) = { 0 if y\_i w^T x\_i >= 0, -y\_i x\_i else

Objective: w\* := arg min\_w l(w; D) + lambda C(w)

Support Vector Machine (SVM) (Hinge & C(w) = ||w||\_2^2)

l\_H(w; y\_i, x\_i) := max(0, 1 - y\_i w^T x\_i)

grad l\_H(w; y\_i, x\_i) = { 0 if y\_i w^T x\_i >= 1, -y\_i x\_i else

L1-SVM (Hinge loss & C(w) = ||w||\_1)

Confusion matrix table with TP, FP, FN, TN. Includes formulas for Accuracy, Precision, Recall/TPR, FPR, F1, and ROC Curve.

k-nearest neighbors (nonlinear) no weights & training, class. during test time. High dim works bad (distance up). Large n to perform well but O(nd) -> O(n^rho), rho < 1 by allowing error probability.

Decision boundary curvilinear. missing data calculate dist. to all neighbors and classified based on values from closest k neighbors. Choose k via CV. Small k -> down bias, up variance (overfit)

y = sign(sum\_{i=1}^n y\_i [x\_i among k nearest neighbors of x])

Decision Tree (nonlinear) no. level/depth chosen unwisely -> small leaf nodes, overfit to noise. Greedy with bad top nodes

Kernel Methods k(x, x') = phi(x)^T phi(x')

K := { k(x\_1, x\_1) ... k(x\_1, x\_n), ... k(x\_n, x\_1) ... k(x\_n, x\_n) } p.s.d <=> for all x in R^n: x^T K x >= 0

Composition k(x, x') as @k\_1 + k\_2 @k\_1 \* k\_2 @c \* k\_1 for c > 0 or f(k\_1), with f = exp or polynomial with pos coeff.

Important Kernels

- k(x, y) = (x^T y + C)^n, with C >= 0, n >= 1
k(x, y) = exp(-||x-y||\_2^2 / (2 sigma^2)), sigma > 0 (Gaussian)
k(x, y) = exp(-alpha ||x - y||\_p), alpha > 0 (Lapl. p = 2, Abel p = 1)

Feature Map may be asymmetric, used to find correct mapping given data. If induced by a valid kernel, may be infinite-dimensional

Neural Network 1 z := w' x 2 v = phi(z) 3 f = w v (Forward)

- linear NN with id more layers => complexity up

Back Propagation

grad\_w l(w\_i, x, y) = [ grad\_w l, grad\_w l ]

Diagram showing flow from x to w to v to f, and corresponding chain rule equations for gradients.

Table with columns Activation, phi(z), phi'(z). Rows include Sigmoid, Tanh, and ReLU (nonlinear).

Vanishing Gradient phi'(z) approx 0 Keep variance of weights approximately constant across layers to avoid vanishing and exploding gradients and network activations (most-prone to VG sigmoid/tanh > ReLU)

Zero-centered E[f(X)] = 0 (id, tanh symmetric around zero)

CNN: For each dimension, the resulting dimension:

l = (padding + filter - stride) / stride + 1

- Number of trainable parameters (default) a filter has the same number of channels as the input

#trainable = #filter \* sum\_i #filter\_dim\_i \* #input channels

Unsupervised Learning

k-Means: CV works badly, low loss if centers close to test set

L(mu) = L(mu\_1, ..., mu\_k) := sum\_{i=1}^n min\_{j in [k]} ||x\_i - mu\_j||\_2^2
mu-hat = arg min\_mu L(mu) non-convex, NP-hard

Lloyd's Heuristics: ~ exponential, in practice not that bad

Initialize cluster center mu^{(0)} := [ mu\_1^{(0)}, ..., mu\_k^{(0)} ]

While not converged, each iteration O(nkd)

- assign points to closest center z\_i^{(t)} -> arg min\_j ||x\_i - mu\_j^{(t-1)}||\_2^2
update centers to mean of each cluster mu\_j^{(t)} -> 1/n\_j sum\_{i: z\_i^{(t)} = j} x\_i

as Hard-EM assume identical, spherical covariance matrices, uniform weights over mixture components. as soft-EM same assumption, with additionally variances -> 0

k-Means++: start with random point as center and add centers randomly, propor. to the squared distance to closest center. Opt expected cost within O(log k) (L(mu) <= O(log k) min\_mu L(mu))

Dimensional Reduction

**PCA Problem:**  $W \in \mathbb{R}^{d \times k}$  orthogonal and  $x_i \in \mathbb{R}^d, z_i \in \mathbb{R}^k$

$$(W, z_1, \dots, z_n) = \arg \min_{W^T W = I_k} \sum_{i=1}^n \|W W^T x_i - x_i\|_2^2$$

$$= \arg \min \sum_{i=1}^n \|W z_i - x_i\|_2^2$$

**Special Case with Closed-form Solution** ( $k > 1$ )

- centered data  $\mu = \frac{1}{n} \sum_{i=1}^n x_i \stackrel{!}{=} 0$  and  $\Sigma = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$

**PCA Solution** non-convex, both  $w$  and  $-w$  are optimal solution

$$W := \underbrace{(v_1 \dots | v_k)}_{1 \leq k \leq d} \text{ and } z_i = W^T x_i \text{ hereby } \Sigma = \sum_{i=1}^d \lambda_i v_i v_i^T$$

**Kernel PCA**

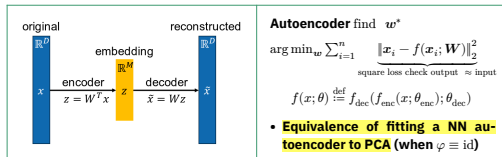
- $z_i := \sum_{j=1}^n \alpha_j^{(i)} k(x, x_j)$
- $\alpha_j^{(i)}$  is the  $j$ -th component of vector  $\alpha^{(i)} \stackrel{\text{def}}{=} \frac{1}{\sqrt{\lambda_i}} v_i$
- $K = \sum_{i=1}^n \lambda_i v_i v_i^T$

**Projection Matrix**  $\underbrace{X}_{\in \mathbb{R}^{m \times n}} \underbrace{X^T}_{\in \mathbb{R}^{n \times m}} \underbrace{b}_{\in \mathbb{R}^m} = X \begin{pmatrix} x_1^T b \\ \vdots \\ x_n^T b \end{pmatrix} = \sum_{i=1}^n x_i x_i^T b$

- Orthogonal square matrix fulfills  $P^2 = P$

**Find projected point** given eigen vector  $v_1$

$$x^{\text{PCA}} := \frac{v_1 v_1^T}{\|v_1\|_2} x = \frac{v_1^T x}{\|v_1\|_2} v_1$$



**Generative Models**  $p(x, y)$  can be more powerful (detect outliers, missing values) with met assumptions, typically less robust against outliers

**Discriminative Models**  $p(y|x)$  detect outliers, but more robust

**GANS** finds saddle point instead of local minimum

**Probabilistic Modeling: choose distribution family  $\mathcal{P}$**

- conditional  $\mathbb{P}[A, B] = \mathbb{P}[A|B] \cdot \mathbb{P}[B]$ 
  - $p(x_i, y_i; \theta) = p(y_i | x_i; \theta) p(x_i)$
- Independent  $\mathbb{P}[A, B] = \mathbb{P}[A] \mathbb{P}[B]$
- Bayes' rule

$$\mathbb{P}[B_i | A] = \frac{\mathbb{P}[A|B_i] \mathbb{P}[B_i]}{\mathbb{P}[A]} = \frac{\mathbb{P}[A|B_i] \mathbb{P}[B_i]}{\sum_{j=1}^n \mathbb{P}[A|B_j] \mathbb{P}[B_j]}$$

$$p_{X|Y}(x, y) = \frac{p_{Y|X}(y|x) p_X(x)}{\int p_{Y|X}(y|x') p_X(x') dx'}$$

$$\mathbb{E}[X] := \begin{cases} \sum_{x \in E} x \cdot \mathbb{P}[X = x] & \text{if discrete } X : \Omega \rightarrow E \\ \int_{-\infty}^{\infty} x \cdot f(x) dx & \text{if cont. } X : \Omega \rightarrow \mathbb{R} \end{cases}$$

$$\sigma^2 = \text{Var}[X] = \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$$

$X \sim$	$p(x) = \mathbb{P}[X = x]$	$\mathbb{E}[X]$	$\text{Var}[X]$
Ber( $p$ )	$p^x (1-p)^{1-x}$	$p$	$p(1-p)$
Bin( $n, p$ )	$\binom{n}{k} p^k (1-p)^{n-k}$	$np$	$np(1-p)$
Poisson( $\lambda$ )	$e^{-\lambda} \frac{\lambda^k}{k!}$	$\lambda$	$\lambda$
Geom( $p$ )	$p(1-p)^{k-1}$	$\frac{1}{p}$	$\frac{1-p}{p^2}$
$\mathcal{U}([a, b])$	$\begin{cases} \frac{1}{b-a} & \text{if } x \in [a, b] \\ 0 & \text{else} \end{cases}$	$\frac{a+b}{2}$	$\frac{1}{12}(b-a)^2$
Exp( $\lambda$ )	$\begin{cases} \lambda e^{-\lambda x} & \text{if } x \geq 0 \\ 0 & \text{else} \end{cases}$	$\frac{1}{\lambda}$	$\frac{1}{\lambda^2}$
$\mathcal{N}(\mu, \sigma^2)$	$\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$	$\mu$	$\sigma^2$

**Multivariate Gaussian**, empirically,  $\hat{\Sigma} \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n x_i x_i^T$

$$\frac{1}{2\pi \sqrt{|\Sigma|}} \cdot \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

**Statistical Inference: Estimating  $\mathbb{P}_{X,Y}^\theta$**

- we used **discriminative model**, where  $\theta$  only parameterizes  $\hat{\mathbb{P}}_{Y|X}^\theta \approx \mathbb{P}_{Y|X}^\theta$

**KL-Divergence**  $D_{KL}(P||Q) = \mathbb{E}_P[\log(\frac{p(x)}{q(x)})]$  (asymmetric)

**Likelihood function** of  $x \in \mathbb{R}^n$

$$L(\theta) = \mathbb{P}[X_1 = x_1, \dots, X_n = x_n] \stackrel{\text{i.i.d.}}{=} p_\theta(x_1) \dots p_\theta(x_n)$$

- e.g.  $p_{Y|X}(D; \theta) = \prod_i p_{Y_i|X_i}(y_i | x_i; \theta)$
- $p(x_{1:n}, z_{1:n} | \theta) = \prod_{i=1}^n p(x_i, z_i | \theta) = \prod_{i=1}^n p(x_i | z_i) \cdot p(z_i)$

**Frequentist Paradigm:** experiments-based only. Precision of estimator unknown if experiment only performed once!

$$\hat{\theta}_{\text{MLE}} = \arg \max_{\theta \in \Theta} p(\mathcal{D}; \theta) \stackrel{\text{i.i.d.}}{=} \arg \max_{\theta \in \Theta} \prod_{i=1}^n p(x_i, y_i; \theta)$$

$$= \arg \min_{\theta \in \Theta} \sum_{i=1}^n -\log \frac{p(x_i, y_i; \theta)}{p(y_i | x_i; \theta) p(x_i)}$$

$$= \arg \min_{\theta \in \Theta} \sum_{i=1}^n -\log p(y_i | x_i; \theta)$$

**Bayesian Paradigm: with "prior belief"**

$$\frac{p(\theta | \mathcal{D})}{\text{posterior belief}} := \frac{p(\mathcal{D} | \theta)}{\text{update}} \frac{p(\theta)}{\text{prior belief (}\exists \text{bias)}}$$

- $\odot p(\mathcal{D} | \theta) \stackrel{\text{def}}{=} \prod_i p(X_i = \dots | \theta)$
- $\odot p(\mathcal{D}) \stackrel{\text{def}}{=} \int p(\mathcal{D} | \theta) p(\theta) d\theta$  (e.g.  $\int_0^1 d\theta$  for Bernoulli)
- posterior distribution  $\mathbb{P}_{\theta|\mathcal{D}}$  computationally hard to compute!

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta \in \Theta} \frac{p(\theta | \mathcal{D})}{p(\mathcal{D}) p(\theta)} \stackrel{\text{i.i.d.}}{=} \arg \max_{\theta \in \Theta} \left( \prod_{i=1}^n p(x_i, y_i | \theta) \right) \cdot p(\theta)$$

$$= \arg \min_{\theta \in \Theta} \sum_{i=1}^n -\log \frac{p(x_i, y_i | \theta)}{p(y_i | x_i, \theta) p(x_i | \theta)}$$

$$= \arg \min_{\theta \in \Theta} \sum_{i=1}^n -\log p(y_i | x_i, \theta) - \log(p(\theta))$$

**Equivalence:** MAP with uniform prior coincides with MLE

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta \in \Theta} p(\mathcal{D} | \theta) p(\theta) \stackrel{!}{=} \arg \max_{\theta \in \Theta} p(\mathcal{D}; \theta) = \hat{\theta}_{\text{MLE}}$$

**Bayes Optimal Predictor: optimal when knowing  $\mathbb{P}_{Y|X}$**

$$f^*_{\hat{f}}(x) \stackrel{\text{def}}{=} \arg \min_{a \in \mathcal{Y}} \frac{\mathbb{E}[\ell(a, Y) | X = x]}{\mathbb{E}}$$

$$= \arg \min_{a \in \mathcal{Y}} \int p(y|x) \cdot \ell(a, y) dy$$

**In practice** when  $\mathbb{P}_{Y|X}$  unknown due to finite dataset, one estimate it by replacing with  $\hat{f}$ , and  $\hat{p}(y|x)$  is obtained from  $\hat{\mathbb{P}}_{Y|X}$  using prob. modeling

**Gaussian Mixture Model**

Weights

$$\sum_{i=1}^K w_i \stackrel{!}{=} 1$$

Overall probability

$$p(x) = \sum_{j=1}^K w_j \overset{\text{probability of } j\text{-th Gaussian}}{\mathcal{N}(x; \mu_j, \Sigma_j)}$$

**Hard-EM** Initialize the parameter  $\theta^{(0)}$ . For  $t = 1, 2, \dots$

- E-Step:** predict most likely class for each data points

$$z_i^{(t)} := \arg \max_z P(z | \theta^{(t-1)}) P(x_i | z, \theta^{(t-1)})$$

- Now data is complete!  $D^{(t)} = \{(x_1, z_1^{(t)}), \dots, (x_n, z_n^{(t)})\}$
- M-Step:** compute closed-form MLE as for Gauss. Bayes classifier

$$\theta^{(t)} = \arg \max_{\theta} P[D^{(t)} | \theta]$$

- i.e.  $\mu_j^{(t)} = \frac{1}{n_j} \sum_{i: z_i = j} x_i$

**Soft-EM (EM):** better deals with overlapping clusters

- Initialize parameters  $\mu^{(0)}, \Sigma^{(0)}, w^{(0)}$ . While not converged:

- E-Step:** prob. sample  $i$ -th belongs to Gaussian  $k$ -th

$$\gamma_k^{(t)}(x_i) = \frac{w_k^{(t-1)} \mathcal{N}(x_i | \mu_k^{(t-1)}, \Sigma_k^{(t-1)})}{\sum_{i=1}^K w_j^{(t-1)} \mathcal{N}(x_i | \mu_j^{(t-1)}, \Sigma_j^{(t-1)})}$$

**M-Step:** fit clusters to weight points. (closed form MLE sol)

<b>k-th weight: ave. prob. that a point belongs to Gaussian</b>
$w_j^{(t)} = \frac{1}{N} \sum_{i=1}^N \gamma_j^{(t)}(x_i)$
<b>k-th mean: weighted average of all points</b>
$\mu_j^{(t)} = \frac{\sum_{i=1}^N \gamma_j^{(t)}(x_i) x_i}{\sum_{i=1}^N \gamma_j^{(t)}(x_i)}$
<b>k-th variance: weighted variance <math>\sigma^2</math> of all points</b>
$\Sigma_j^{(t)} = \frac{\sum_{i=1}^N \gamma_j^{(t)}(x_i) (x_i - \mu_j^{(t)})^2}{\sum_{i=1}^N \gamma_j^{(t)}(x_i)}$
<ul style="list-style-type: none"> <li>to avoid degeneracy, one can add <math>\nu^2 I</math> to the diagonal of MLE in update, equivalent to placing a (conjugate) Wishart prior on the covariance matrix</li> </ul>

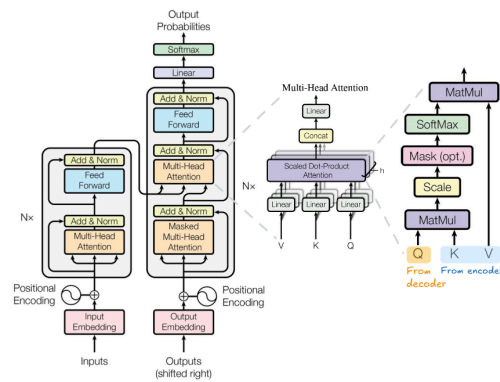
**LLM**

**RNN:** Recall **feedforward NN** neglecting

- order context  $\Rightarrow$  cannot capture sequential data
- RNN** has recurrent connection,  $\text{input}^{(t+1)} := \text{output}^{(t)}$  (embedding memory, eg. long-short-term memory **LSTM**)
- Sequence-to-Sequence Model Disadvantage:** long path of info flow, hard to capture long range dependencies (vanishing gradient)

**Transformer** process in **parallel** where everything can be passed at once, efficiency allows to train more data (**Note:** but no direct effects on training data set)

- Architecture:** also with Residual skip connections; **multi-headed blocks** for diff. meanings



**Matrix Calculation of Self-Attention** With word embedding matrix  $X$ , apply **MatMul** with trained weights  $(W^Q, W^K, W^V)$  and get  $(Q, K, V)$ . Then obtain matrix product  $Z$  with:  $\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) \cdot V = Z$

Layer Type	Complexity per Layer	Sequential Operations	Maximum Path Length
Self-Attention	$O(n^2 \cdot d)$	$O(1)$	$O(1)$
Recurrent	$O(n \cdot d^2)$	$O(n)$	$O(n)$
Convolutional	$O(k \cdot n \cdot d^2)$	$O(1)$	$O(\log_k(n))$
Self-Attention (restricted)	$O(r \cdot n \cdot d)$	$O(1)$	$O(n/r)$

**Positional Encoding:** attention neglected orders! Encode with trig. functions; each position/index is mapped to a vector, output of the layer is a matrix with each row as an encoded object

**Fine-Tuning:** traditionally requires gradient updates, **Methods with no grad updates:** Zero-shot predicts the answer given only a natural language description of the task. One-shot a single example of the task provided Few-shot a few examples provided

1	Translate English to French:	← task description
2	sea otter => loutre de mer	← example
3	cheese =>	← prompt