

- Class Notes - Intro

Like DevOps, automation & workflow in ML \Rightarrow MLOps but it's
MLOps combines ML, DevOps & Data Engineering late + data

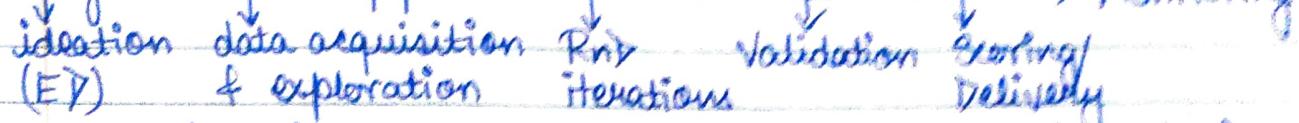
Kubeflow pipelines - visualize ML pipelines

Tensorflow transform - check consistency

Automate unit tests & integration test \rightarrow TF data validation

- Descriptive • Predictive • Prescriptive • Retrospective

Key items in ML pipeline



• Vertex AI - Google

New ML-problem \rightarrow Optimized Marketing using RL

Characteristic	NN	SVM	Trees	MARS	k-NN, kernels
Handling mixed data	↓	↓	↑	↑	↓
Missing Values	↓	↓	↑	↑	↑
Outliers	↓	↓	↑	↓	↑
Insensitive to ϕ	↓	↓	↑	↓	↓
Scalability	↓	↓	↑	↑	↓
Irrelevant inputs	↓	↓	↑	↑	↓
Linear comb ⁿ of X	↑	↑	↓	↓	~
Interpretability	↓	↓	~	↑	↓
Predictive power	↑	↑	↓	~	↑

MARS \rightarrow auto detects non-linearity & interaction & non-parametric

• Neural net issues \rightarrow ↑ data, ↑ hyp. tuning \Rightarrow ↑ compute & stats inc.

• Correlation Causation

• IHME \rightarrow outside 95% CI, 70 percent times \Rightarrow ↑ noise \Rightarrow ↓ SNR.

Curvefit - Non linear mixed effects \rightarrow fixed ① consideration

* For fixed effect Fisher Information random matrix used to get asymptotic approx of uncertainty.

can understand Fixed vs Random from Simpson's paradox case.

$$\text{Gaussian density } N(\alpha | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\cdot z^2}$$

$$E(x) = \mu \quad E(x^2) = \mu^2 + \sigma^2$$

$$\text{Var}(x) = \sigma^2$$

Likelihood = \prod density of i instances (under independence assumption)

overall

Maxim Likelihood Estimation -

$$\hat{\theta}_{MLE} = \underset{\theta}{\operatorname{arg\ max}} p(X|\theta)$$

Distribution

Exponential

$$\begin{aligned} \theta > 0 & \text{ pdf: } \lambda e^{-\lambda x} \\ \theta > 0 & \text{ cdf: } 1 - e^{-\lambda x} \\ & \text{ median: } \ln(n) \\ & \text{ mean: } (1/\lambda) \\ & \text{ var: } (1/\lambda^2) \end{aligned}$$

$$\begin{aligned} \text{free metric pdf: } & (1-p)^{k-1} p \\ \text{cdf: } & 1 - (1-p)^k \\ & \text{ median: } -1/\log_2(1-p) \\ & \text{ mean: } (1/p) \end{aligned}$$

$$\text{Normal pdf: } \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}z^2}$$

$$\mu \sim \sigma^2 \text{ pdf: } \frac{1}{\sqrt{2\pi}\sigma} \left(1 + \frac{(x-\mu)^2}{2\sigma^2} \right)^{-\frac{1}{2}}$$

where erf(x) = $\frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt$

$$\text{Binomial pdf: } n! x^k (1-p)^{n-k}$$

Poisson $\lambda^x e^{-\lambda} / x!$

$$\hat{x}_{MLE} = \frac{\sum x_i}{n}$$

Log Likelihood

$$\begin{aligned} -n\lambda + \sum_{i=1}^n \lambda x_i &= d[\ln L(\lambda)] \\ \Rightarrow \max_{\lambda} \lambda \sum_{i=1}^n x_i &= \lambda \bar{x} \end{aligned}$$

But \bar{x} is unbiased est of p

$$\begin{aligned} \text{free metric pdf: } & (1-p)^{k-1} p \\ \text{cdf: } & 1 - (1-p)^k \\ & \text{ median: } -1/\log_2(1-p) \\ & \text{ mean: } (1/p) \end{aligned}$$

$$\begin{aligned} \ln L(\lambda) &= -\frac{n}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma^2} \sum (x_i - \mu)^2 \\ \text{wrt } \mu &\rightarrow \frac{1}{\sigma^2} \sum (x_i - \mu) \Rightarrow \boxed{\bar{x}} \text{ for zero} \\ \text{wrt } \sigma^2 &= -\frac{n}{2} + \frac{1}{2\sigma^2} \sum (x_i - \mu)^2 \Rightarrow \text{for zero} \\ \sigma^2 &= \frac{1}{n} \sum (x_i - \bar{x})^2 \min \text{var (unbiased)} \end{aligned}$$

$$B(n, p) \rightarrow N(np, npq) \text{ for } \uparrow n \text{ if } \uparrow p$$

$$B(n, p) \rightarrow Pois(np) \text{ for } \uparrow n \text{ if } \uparrow p$$

Comments

pdf of "time last" events in Poisson point process
 → event count of independent events at constant avg. rate.
 Special case of gamma (continuous analogue of geometric)
 Approx. bias corr = $\frac{n-2}{(\bar{x})}$ for MLE

trials until success: discrete
 Assump: Independence, binary, constant p
 to model

\Rightarrow Inflections at $H \pm \sigma$
 Cauchy($\text{df}=2$) → Student's t → Normal ($\text{df}=\infty$)

Minimizing MSE on train data in regression
 using standard MLE solution of assumed
 linear model

$$B(n, p) \rightarrow N(np, npq) \text{ for } \uparrow n \text{ if } \uparrow p$$

$$B(n, p) \rightarrow Pois(np) \text{ for } \uparrow n \text{ if } \uparrow p$$

(Bayesian) Conjugate Priors

* Generally 1 hyp. param ↑

Posterior = Likelihood * Prior
 (Beta) (Binomial) (Beta) \Rightarrow prior is CP of likelihood

Likelihood Conjugate Prior Posterior Hyp. Interpretation

Bernoulli (p) Beta (α, β) $\alpha + \bar{x} \cdot n, \beta + n - \bar{x} \cdot n$, α success, β failures

Binomial (p) Beta (α, β) $\alpha + \bar{x} \cdot n, \beta + \Sigma n - \bar{x} \cdot n$ " "

Poisson (λ) Gamma (α, β) $\alpha + \bar{x} \cdot n, \beta + n$ total occurrence in binomial

Geometric (p) Beta (α, β) $\alpha + n, \beta + \bar{x} \cdot n$ α experiment, β failures

Multinom

Dirichlet

Normal (μ, σ^2) Normal complicated

Uniform

Pareto

Pareto

Gamma

Exponential

Gamma

$\alpha + n, \beta + \bar{x} \cdot n$

Multivariate Gaussian Distribution:

$$N(\boldsymbol{\alpha} | \boldsymbol{\mu}_{\Sigma}) = \frac{1}{(2\pi)^D \sqrt{\det \boldsymbol{\Sigma}}} e^{-\frac{1}{2} (\boldsymbol{\alpha} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\alpha} - \boldsymbol{\mu})}$$

* Marginal +
Conditional prob. of
normal are normal
themselves

- ML Short new

$$\text{Anova} \rightarrow F_{\text{stat}} = \frac{\text{MSR}}{\text{MSE}} \quad t_{\text{stat}} = \left(\frac{\beta - \bar{x}_{\text{test val}}}{\text{SE}} \right) \quad \text{Adj } R^2 : \text{penalize } n \uparrow$$

OLS \rightarrow sales \sim price₁ + price₂ + ϵ_1 Nested impl of plm in R
Check?

Fixed effects \rightarrow sales \sim p₁ + p₂ + city fixed effects + ϵ_2

Random effects \rightarrow sales \sim p₁ + p₂ + xfe + time fixed effects + ϵ_3

can model \uparrow cov. as

RE $\sim N(\mu, \sigma^2)$ agnostic

within

city

individual

intercity +

within city

(idiosync error)

$$\hat{\sigma}_{b_j}^2 = \frac{s^2}{(N-1)(\text{var in } x_j \text{ not ass. with other } x_i)}$$

- Succinct models with superior features are accurate, robust

* Function Approximation - Regression

$y(x, w) = w_0 + w_1 \phi_1(x) + \dots + w_n \phi_n(x)$ to predict t (true target)
 prediction weights for dependent variable \downarrow independent basis funcⁿ variables

$E(t)$ given ϕ_i is linear in $\phi \Rightarrow E(t|\phi) = \text{Linear Combos } (\phi_i)$

- parametric form

Assumptions

- $E(t|\phi)$ linear in ϕ - Errors random $N(0, \sigma^2)$ - No multicollinearity
- Exogeneity in ϕ - Homoskedasticity - No serial autocorrelation

• Check for Simpson's paradox especially in panel data

Test for Normality - Kolmogorov-Smirnov (CDF) Shapiro-Wilk for better

Test for multicollinearity - Breusch Pagan (test st. is Chi-sq distributed)

Test for autocorrelation - Durbin Watson (DB only checks with lag 1)

Test for (error) endogeneity - Hausman

MLE provides error metric as squared loss for best estimator

* MVUE ~ non-linear BLUE

* Exact closed form soln \Rightarrow iterative approach not needed

$$\vec{w}^* = (\phi^T \phi)^{-1} \cdot \phi^T \vec{t}$$

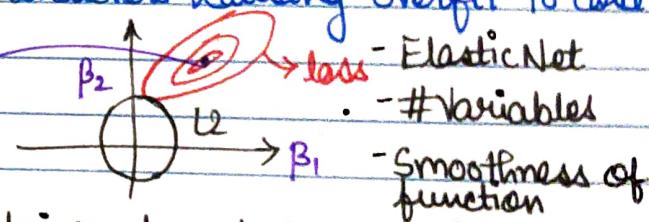
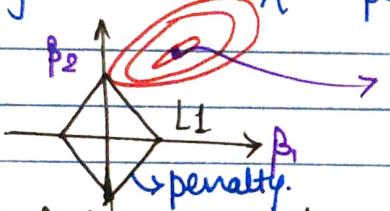
If collinear $\Rightarrow \det(\phi^T \phi) \rightarrow 0 \Rightarrow$ inverse

Moore Penrose Inverse becomes large \Rightarrow parameter estimates have higher uncertainty (No effect on R^2 cause info. is there)

SEs are defined as -

$$\delta_{bj} = \frac{s^2}{(N-1) \text{ (Variation in } x_j \text{ not associated with other } x_i\text{'s)}}$$

Regularization: Add penalty to whatever is causing overfit to cost f



As λ increases, chance of betas being shrunk to zero increases.

- Intercept not included case - ① No need, intercept should be close to $E(y_i)$ ② If regularized prediction of model with $y_i + c \sim \forall x_i$ would not necessarily increase by c .

• Standardize all x_i , else penalty could be overpowered.

- Deviance vs $-\log \lambda$ graphs with variance error bars + λ_{\min} + λ_{\max}
- called weight decay in NN - shrinkage in stats

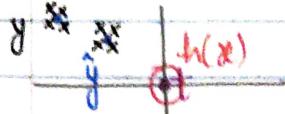
$$\vec{w}^* = (\lambda I + \phi^T \phi)^{-1} \cdot \phi^T \vec{t} \Rightarrow \text{adds numerical stability}$$

Model evaluation $\rightarrow f(\text{train}, \text{validation}) \Rightarrow$ collection of solutions \Rightarrow CIs can be
 • Bias Variance Trade-off \hookrightarrow not in Bayesian \Rightarrow No CIs \hookrightarrow no probabilistic distribution

Say $t = h(x) + \epsilon$ where $\epsilon = N(0, \sigma^2) \Rightarrow h(x)$ is best possible function in \mathcal{H}

Our prediction $= y(x) \Rightarrow$ Expected Loss $= \int (y(x) - h(x))^2 p(x) dx$ (bias² + variance)
 for dataset \rightarrow we have $\left[(h(x) - E[y(x)])^2 + E[(y(x)) - E(y(x))]^2 \right] + \int \int (h(x) - t)^2 p(x, t) dx dt$ (irreducible error)

- * t not function of $x \Rightarrow$ double integral
- * Given $h(x)$ is unique, bias² is not
- * Given multiple models \Rightarrow variance is double expectation enclosed in another expectation

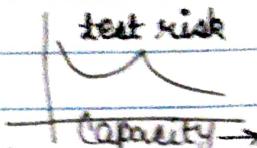


High bias, low variance



High variance, low bias

\rightarrow property of model ω



Capacity \rightarrow

* Regularization due to over parameterization in neural networks

Explaining Bias-Variance trade-off \rightarrow # bins in histogram.

Incremental Forward Stagewise Regression -

Residual $r \leftarrow y$, and $\forall i, \beta_i = 0$; $\forall x_i$ we have $\sum x_i = 0$ + $\sum x_i^2 = 1$ (unit norm)

$[\beta_i \leftarrow \beta_i + \epsilon \cdot \text{sign}(x_i, r)]$ $\max_{\text{correlation}}$ \Rightarrow until r uncorrelated with all x .

$f[r \leftarrow r - \epsilon \cdot \text{sign}(x_i, r) \cdot x_i]$ Q. shouldn't this be repeated too many times, but seems $O(n)$?

Other Regressions -

\rightarrow Partial Least Squares

Logistic, Polynomial,

Quantile (Median = least abs. dev.)

Elasticnet

\rightarrow Principal Components Regression \rightarrow Support Vector Regression

\rightarrow Poisson Regression vs Negative Binomial vs Quasi poisson regression
 $(\mu = \sigma^2 = \lambda)$ $(\mu \neq \sigma^2)$ but Quadratic $\sigma^2 = \text{Quad}(\mu)$ $\sigma^2 = \text{Linear}(\mu)$

\rightarrow Cox Regression (Coxph \rightarrow proportional hazards generally never satisfied)

\rightarrow Tobit Regression

\rightarrow Least Angle Regression (Efficient Lasso)

\rightarrow Panel Regression

Succinct models with 'right' predictors are superior.

- Transformation types → Log (in case of outliers in X or subsequently in residuals) → Standardization ($\bar{x} - \mu$) → Minmax (shift to 0,1 range) → Box Cox transformation → $\log x + \sqrt{\lambda=0}$ {find closest to gaussian} → Yeo Johnson (similar but separate treatment for negatives) → Rank Gauss (values → Rank → Gaussian) ** → Freq to Rank encoding; Target encoding; Binary encodings

Almost linear methods -

- Additive models - GAMs
- Piecewise linear - MARS
- Locally linear - LOESS

* Universal Approximator -

✓ continuous function, member of the family exactly matches it.

Often $E(\theta)$ is non-linear, not a quadratic f often non-convex → SGD

Incrementally 'surf' down the cost function in weight space

$$w^{(t+1)} = w^{(t)} - \eta \cdot \nabla E \quad \text{Adaptive eta - } \textcircled{1} \text{ Nelder-Mead}$$

Scaling Nelder Mead to n dimensions - use $(n+1)$ points convex poly. Contraction vs Expansion

Can become divergent in any edge case?

Gradient Descent → Elegant move towards nearest minima (can take time) will oscillate if eta set too high. Line Search somehow

adapt eta so that loss is always decreasing **

* SGD faces issue with Rosenbrock: $(1-x)^2 + 100(y-x^2)^2$

(getting stuck for few inits & slow) → steep & variable change in dimensions

Conjugate Gradient - Take into account previous search direcⁿ (momentum?)

$$S_n = \nabla X_n + \beta S_{n-1} \text{ where } \beta = \frac{\nabla X_n^T \cdot \nabla X_n - \nabla X_{n-1}^T \cdot \nabla X_{n-1}}{\nabla X_{n-1}^T \cdot \nabla X_{n-1}}$$

using search direction f component from line search of previous search direction we have next x location as $X_t = x^* + \alpha S_n$

* Helps in cases like Rosenbrock. Is momentum about acceleration f avoids the zigzag problem conjugate gradient about direction?

PCA = MDS (Euclidean Dist)

↓ eigendecompose $\xrightarrow{\text{Dist.}} \tilde{X} \tilde{X}^T \xrightarrow{\text{Covariance mat.}} (\tilde{X} \tilde{X}^T) \cdot \tilde{V} = \tilde{\Lambda} \tilde{V}$ MDS = Principal Coordinates Analysis (PCoA)

loss function strain - difference between

$(\tilde{X}^T \tilde{X}) \cdot \tilde{V} = \tilde{\Lambda} \tilde{V} \xrightarrow{\text{actual f wanted distances}} (\tilde{X} \tilde{X}^T) \cdot \tilde{V} = \tilde{\Lambda} \tilde{V}$ **

Both are giving same $\tilde{\Lambda}, \tilde{V}$ translates to eigendecomposition in case of Euclidean

conjugate gradient → faster than line search, more smoothed as well.

SGD → better in online sense → streaming data → non-stationary

Ideally decrease η with time to allow convergence (a/t)

Batch → #epochs* (one pass over all data where 1 batch taken in 1 go)

SGD → faster, better for large data, ↓ prone to local minima, stationary is fine

Given $G \triangleright$ is slow \Rightarrow leverage 2nd order \Rightarrow Newton's method \Rightarrow Quadratic approximation of Taylor in 1d
In multi-dimension -
 $\hat{x}_{n+1} = \hat{x}_n - \nabla(f(\hat{x})) \cdot H(f(\hat{x}))^{-1}$
could move Jacobian (gradients) Hessian (curvatures)
towards global wrt \hat{x} wrt \hat{x} } Difficult to compute the Hessian & inverse issues.

minima given assumption holds true \Rightarrow Big steps in low curvature & vice versa.

Better to have line search SGD + leverage previous directions \Rightarrow Conjugate grad

* Neural networks: Adaline vs Perceptron } Both are classifiers
weights = # pairs $\rightarrow t - f(x) \rightarrow t - \text{step}(f(x))$ } slight difference
biases = # perceptrons (# linear combinations)

Multi-layered perceptron \rightarrow if output is linear & hidden are non-linear
 \Rightarrow Universal Approximator

total weights = (Input+1)Hidden₁ + (Hidden₁+1)Hidden₂ + (Hidden₂+1)Output

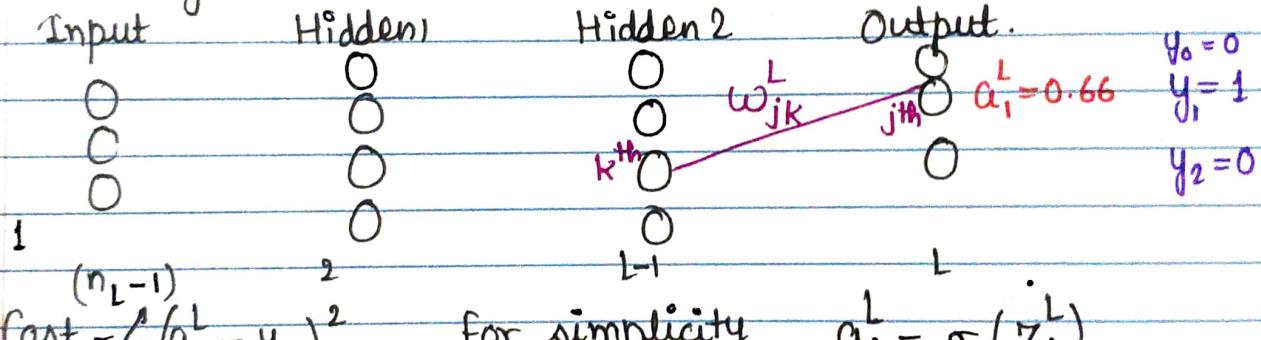
why logistic function for sigmoid curve? $\sigma(x) = \frac{1}{1+e^{-x}}$ $\sigma'(x) = \sigma(x)(1-\sigma(x))$

Forward pass - Given input calculate output layer by layer

Compute error at final output layer

Backward pass - Compute gradient of weights wrt loss, layer by layer

Update all weights.



$$\text{Cost} = \sum_{j=0}^{(n_L-1)} (a_j^L - y_j)^2$$

(or average)

for simplicity

$$z_j^L = \sum_{k=0}^{(n_{L-1}-1)} w_{jk}^L \cdot a_k^{L-1} + b_j^L$$

$$a_j^L = \sigma(z_j^L)$$

$$w_{jk}^L \cdot a_k^{L-1} + b_j^L$$

\Rightarrow The chain rule becomes -

$$\frac{\partial \text{Cost}}{\partial w_{jk}^L} = \frac{\partial \text{Cost}}{\partial a_j^L} \cdot \frac{\partial a_j^L}{\partial z_j^L} \cdot \frac{\partial z_j^L}{\partial w_{jk}^L}$$

Cost sensitivity wrt all weights - if multiple paths \Rightarrow add those up.

$$2(a^L - y) \cdot \sigma'(z^L) \cdot (a^{L-1})$$

Given this is for two general layers, recursively repeat for all the layers. Use gradient term to update weights.

\Rightarrow Store gradients somewhere to avoid recompute.

SGD → Man walking down the hill - slow, careful & steady
 Momentum → Heavy ball rolling down the same - Newton's 1st law.
 Robustness → RNN in R or SGR → Not a function given 1-M
 Modeling inverse problem? solutions → Model joint probability

* Sequential Learning with SGD → non-stationary OR piecewise model

If you want convergence → η decreases with time Eg. RMSprop.

$\eta(t) = 1/t \Rightarrow$ Robbins-Monroe sequence

RNNs	LSTMs	AEs	VAEs	GRUs
sequence memory via ip \rightarrow h _{t-1} \rightarrow op but no (BPT) forget gate \rightarrow h _{t-1} \rightarrow (long) of (ip+forget) emb \rightarrow op - denoising	memory via memory touch to in update \rightarrow h _{t-1} \rightarrow (long) of (ip+forget) emb	image compression possible	wrote new hand i/p to written i/p to images distrib mapping then sampling	reset & update gates lessen operations ops on hidden states \rightarrow faster

Markov Chain Boltzmann machine
by state transition patterns by reconstr' i/p
matrix

Convolutional NN
basically n-dim NN
filter, stride
padding etc.

GANs*

Graph NNs*

Attention Network*

* Data Preprocessing

Boxplot: Whiskers upto extremes or 1.5 IQR outside IQR whichever nearest

Missing: Check if random/systematic else fill in conditional mean/mode/knn

Outlier detection → Prob. based \Rightarrow Estimate pdf(x) given $x_i; \forall i$

Parzen window \Rightarrow kernel density Estimate \Rightarrow Approx f as mix of continuous(k) distributions centered at x_i with bandwidth h. k is window or kernel

$$\hat{f}_h(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right) \quad K \text{ can be Normal, poission etc.}$$

① Difference between GAMs, KDE and (Gaussian) Mixture models?

Discrimination based: 3 sigma, freq less than 1%, dist $>$ kth neighbour Average distⁿ

* Trimmed mean * Huber loss instead of square loss

* Posthoc \rightarrow Roughly monotonic to monotonic by isotonic regression

PyOD \rightarrow COPOD