

BAYESIAN THEOREM

$$* P(H|E) = \frac{P(E|H) \cdot P(H)}{P(E)} \rightarrow \begin{array}{l} \text{class prior} \\ \text{likelihood} \\ \text{posterior} \end{array}$$

\rightarrow predictor prior

- * It takes the data & calculates probabilities of each E_i separately for all classes.
 - o Normal Mail : $P(E_1), P(E_2) \dots \rightarrow P(E_1|N), P(E_2|N)$
 - o Spam Mail : $P(E_1), P(E_2) \dots \rightarrow P(E_1|S), P(E_2|S)$
- * It calculates all apriori probabilities of $P(N), P(S)$ by simple proportion calculation.
- * Then, calculate $P(N|E)$ etc.. is calculated from the apriori probability and conditional probability.

Naïve: It assumes that evidence splits into parts that are conditionally independent!

Zero Frequency Problem: If one of the attributes has never occurred in one class, its prob $|P(E|S)|$ will be 0. So, it will disrupt.

\Rightarrow ADD or Laplace est. ~ 1 (usually) for each attribute in each class, and calculate probabilities.

- * Generally used in:
 - Real Time Prediction
 - Multiclass Prediction
 - Text classification / sentiment Analysis
 - Recommendation System with collaborative filtering

NOTES OF USAGE

- Don't forget about zero FREQUENCY problem.
- Independence of attributes is a strong assumption, beware!
- It is fast & successful when independence holds.
- Its probability might be meaningless, don't use it.

* If the data is categorical → OK!

If not: **Gaussian Naive Bayes**

→ which assumes variables follow a normal distribution (STRONG ASSUMPTION)

PYTHON

* Sklearn.naive-bayes

- └ GaussianNB() → non-categorical variables
- └ MultinomialNB() → categorical variables
- └ ComplementNB() → imbalanced datasets (adaptation of M. NB)
- └ BernoulliNB() → work with already binary or binarize data → occurrence counted
- └ CategoricalNB() → categorically distributed data

• fit(x-train, y-train) }
• predict(x-test) } for all of them...

Binomial Distribution

* Bernoulli: A random experiment with exactly two possible outcomes \rightarrow Fail / Success

* Binomial: Number of successes for a sequence of Bernoulli trials.

$$* f(x) = \binom{n}{x} p^x (1-p)^{n-x}$$

$$\circ \mu = n \cdot p$$

$$\circ \sigma^2 = n \cdot p \cdot (1-p)$$

\rightarrow Fixed # of trials.

\rightarrow Random # of successes.

Geometric Distribution

\rightarrow Fixed # of successes

\rightarrow Random number of trials

{ For x/ success, how many trials should be done?

$$* f(x) = p (1-p)^{x-1}$$

↳ where x is the number of failures until the 1st success

$$\circ \mu = \frac{1}{p}$$

$$\circ \sigma^2 = \frac{(1-p)}{p^2}$$

\rightarrow Lacks of Memory: So if we start trial number i in the process, it does not differ from starting at 0.

Poisson Distribution

\rightarrow Expresses the probability of a given number of events occurring in a fixed interval of time / space.

$$* f(x) = \frac{e^{-\lambda} \cdot \lambda^x}{x!} \quad (\lambda \text{ is the rate per time / space})$$

$$\circ \mu = \lambda \text{ and } \sigma^2 = \lambda$$

* λ should be events per unit time!

↳ If not, convert it!

Poisson Approximation to Binomial

* Binomial cases when

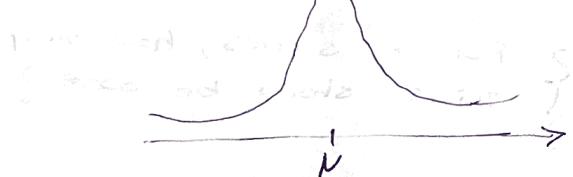
- n is large \uparrow $\left\{ \begin{array}{l} n \approx 10 \\ np < 10 \end{array} \right.$
- p is small \downarrow $\left\{ \begin{array}{l} n \approx 10 \\ np > 10 \end{array} \right.$

THEN:

* Use Poisson with $\lambda = np$

CONTINUOUS

Normal Distribution



$$* f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$\left. \begin{array}{l} \cdot \mu = N \\ \cdot \sigma^2 = \sigma^2 \end{array} \right\} N(\mu, \sigma^2)$$

→ Standard z-table gives cumulative prob. of distribution for certain values in Z .

APPROX. TO BINOMIAL

$$P(X \leq A) \approx P(Z \leq A + 0.5)$$

↳ where $N = np$

$$\sigma^2 = np(1-p)$$

↳ use normal z conversion \rightarrow NORMAL PROBABILITY

APPROX. TO POISSON

$$P(X \leq 950) = \sum_{x=0}^{950} \frac{e^{-1000} \cdot 1000^x}{x!} \quad \underline{\text{WUFF!}}$$

$$\approx P(Z \leq 950.5) \quad \text{and} \quad \begin{aligned} N &= 1000 (\lambda) \\ \sigma^2 &= 1000 (\tau) \end{aligned} \quad \rightarrow \text{NORMAL Z.}$$

Exponential Distribution

- On probabilities of events occurring at a time, or in a range of time.

$$* f(x; \lambda) = \begin{cases} \lambda e^{-\lambda x}, & x \geq 0 \\ 0, & x < 0 \end{cases}$$

$$\begin{aligned} * P(X \leq x) &= 1 - e^{-\lambda x} \\ * \mu &= \frac{1}{\lambda} \\ * \sigma^2 &= \frac{1}{\lambda^2} \end{aligned}$$

- Also has lack of memory property.

Beta Distribution

$$(\alpha, \beta)$$

- A probability distribution on probabilities. $[0, 1]$.

- Binomial models number of successes whereas

Beta models the probability of success

* $\alpha - 1 \rightarrow \# \text{ of successes}$

* $\beta - 1 \rightarrow \# \text{ of failures}$

Gamma Distribution

$$(k, \theta)$$

- Predicts the wait time until the k^{th} event occurs. (whereas exponential predicts the wait time for 1st event)

$k = \# \text{ of events} \text{ (scale)}$

$\theta = \text{shape-rate of events} \text{ (rate)} \quad \left. \begin{array}{l} \text{Positive} \\ \text{constant} \end{array} \right\}$

* $k = 1 \rightarrow \text{EXPONENTIAL}$

* $k = \text{integer} \rightarrow \text{ERLANG}$

STATISTICS

- * Frequentists define probability as the long-term frequency in a repeatable random experiment
- * They do not use probability to quantify incomplete knowledge
- * measure degree of belief in hypothesis

HYPOTHESIS TESTING

We propose:

$$H_0 = \text{Null Hypo}$$

$$H_A = \text{Alternative Hypo}$$

\Rightarrow We try to reject H_0 in the direction of H_A
so that we can draw conclusions.

- One-Sided $P = P(Z > z|H_0)$
- Two-Sided $P = P(|Z| > z|H_0)$

p-value: is the probability that of the values farther to the center than this value.

Ex: A p-value of 0.05 means that there is 5% chance of finding more deviated points.

(GIVEN THAT H_0 !)

- * Whilst doing that, p-value is the probability given that H_0 is true. If $H_0 = \text{N=0}$:

$p = 0.32$ is the test statistic being the closest among 32% of the data if the center is 0.

\hookrightarrow In this case, data suggests that not significantly away from assumed center, not at the tails \rightarrow So, N can be 0. Cannot reject!

ERRORS & MEANINGS

REALITY

	H_0	H_A
Reject H_0	Type I	✓
Don't Reject H_0	✓	Type II

* $\alpha = \text{Significance} = P(\text{Type I})$

* probability we incorrectly reject H_0 . $\{ \sim 0 \vee$

* probability of rejecting H_0 when it is true.

* $\beta = \text{power} = 1 - P(\text{Type II})$

* probability we correctly reject H_0 .

Student t-test

→ When variance is unknown, we have to do an estimation.

$$s^2 = \frac{\sum (x_i - \bar{x})^2}{N-1} \rightarrow N-1 \text{ since } \bar{x} \text{ is an estimation as well.}$$

→ And the t-test becomes:

$$T = \frac{\bar{x} - \mu}{s/\sqrt{n}}$$

Chi-square Distribution

$$\sum \left(\frac{x_i - \mu_i}{\sigma_i} \right)^2 = \chi^2 = \frac{(n-1) \cdot s^2}{\sigma^2}, \quad (n-1 \text{ d.o.f.})$$

STATISTICAL METHODS

SAMPLING

- ↳ Random Sampling
- ↳ Systematic Sampling
- ↳ Stratified Sampling

RESAMPLING

Once we sample the data, we may try to make a statistical inference about the population from a small set of observations.

- ↳ But with one sample we have a single estimate with little idea of the variability in the estimate.
- * There comes resampling!

* NORMALLY SUB-SAMPLES ~ Regular (without replacement)

BOOTSTRAP

- Sampling with replacement
- There are Out-of-Bag (OOB) data for testing / validation purposes.
- MIGHT be used for model development, confidence interval building, estimation etc.
- * (Sample size might be equal to the whole data)

PYTHON

- sklearn.utils
 - ↳ resample (data ,
 - replace = True ,
 - n_samples = ,
 - random_state =)

METRICS

Confusion Matrix

Predicted

		<u>Actual</u>	
		Positive	Negative
Positive	Positive	TP ✓	FP ✗
	Negative	FN ✗	TN ✓

→ Binary classification matrix

Accuracy

$$\frac{TP + TN}{TP + TN + FN + FP}$$

→ How accurate I am?

! Not good in imbalanced data!

Precision

$$\frac{TP}{TP + FP}$$

→ How precise I am when I guess positive?

Recall

$$\frac{TP}{TP + FN}$$

→ How successfully recall the positives when I am guessing?

F1 Score

$$* 2 \cdot \frac{1}{\frac{1}{precision} + \frac{1}{recall}}$$

→ Harmonic mean between precision & recall.

→ How precise your classifier is as well as how robust it is (not missing significant number of instances)

⇒ Here, we are more focused on the positive class where we try to be precise in our predictions and also capture as much as the positive class at the same time.

Log Loss ~ cross-entropy function (used as metric as well)

- works by penalizing the false classifications.
- Works well for multi-class classification.
- Classifier should assign probabilities to each class for all samples.

- * N samples belonging to M classes:

$$\text{LogLoss} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^M y_{ij} * \log(p_{ij})$$

↳ Unlike MSE it gives 0 weight to wrong predictions but heavy weight to correct predictions

↳ When we make it for binary prediction: (m=2)

$$\text{LogLoss} = -\frac{1}{N} \sum_{i=1}^N [y_i \cdot \log(p_i) + (1-y_i) \cdot \log(1-p_i)]$$

- * It actually can be derived from maximum likelihood of bernoulli trial.
- * Since it penalizes distances from the actual values (in prob. of 0-1) it turns out to be a negative value but should be used as a minimization problem

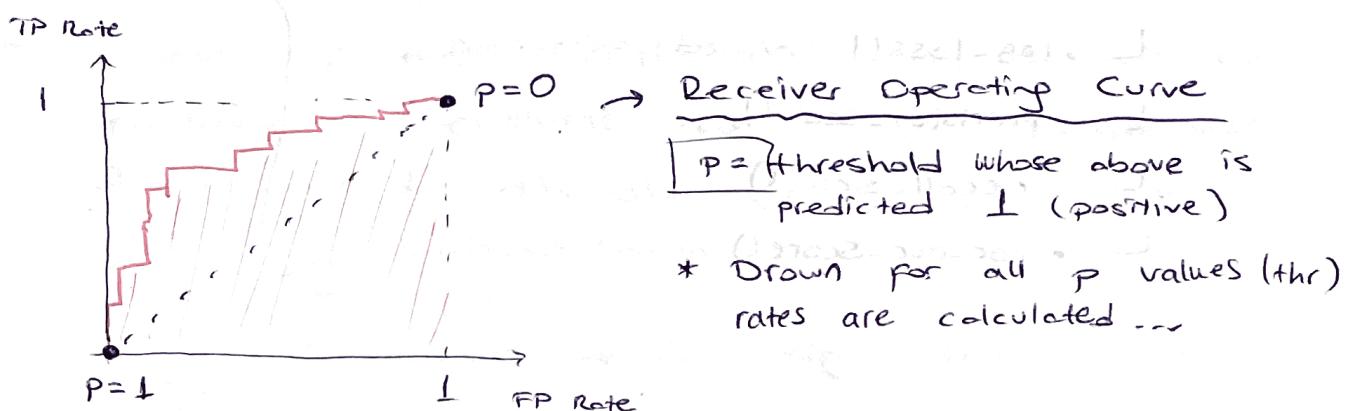
→ THEREFORE NEGATIVE SIGN.

AREA UNDER CURVE

- Known as C-value, as well.
- It is the probability that when one sample from positive class and one from negative class are drawn, the prob. assigned to positive is greater than the negative (ranked higher). → RANKING ORDER...
- TP (Sensitivity) Rate and FP (1-Specificity) rate is compared.

* TP Rate = $\frac{TP}{TP + FN}$ ~ How much of the actual positive I guessed/recalled

* FP Rate = $\frac{FP}{FP + TN}$ ~ How much of actual negative I guessed INCORRECTLY positive.



↳ Area under Curve (AUC) is the metric.

Mean Absolute Error

$$MAE = \frac{1}{N} \sum |y_j - \hat{y}_j|$$

↳ how far the predictions from the actual output.

Mean Squared Error

$$MSE = \frac{1}{N} \sum (y_j - \hat{y}_j)^2$$

R² & Adjusted R² ~ coefficient of determination

* R² is a metric based on how much of the variance in the data is explained by the suggested model.

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}} = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

$$\text{Adj. } R^2 = 1 - \frac{(n-1)}{[n-(k+1)]} \cdot (1-R^2)$$

↳ penalizes addition of useless variables...

PYTHON IMPLEMENTATION

* sklearn.metrics

L. accuracy_score()

L. f1_score()

L. log_loss() ~ need predict-proba

L. precision_score()

L. recall_score()

L. roc_auc_score() ~ need predict-proba

CLASSIFICATION

L. explained_variance_score()

L. max_error()

L. mean_absolute_error()

L. mean_squared_error()

L. median_absolute_error()

L. r2_score()

L. mean_absolute_percentage_error()

REGRESSION

SKLEARN CROSS-VAL

* sklearn.model_selection.cross_val_score

(classifier-object,

X,

y,

cv=k(fold),

scoring='name of scoring method')

W
O

sklearn.metrics.SCORERS is a dict

& .keys() gives allowed metric names.

* sklearn.metrics.classification_report(y_true, y_pred, target_names=C)

↳ gives precision/recall/f1 for each class.

LOSS FUNCTIONS

- * Functions that are defined as the deviation (somehow) of predictions from actual values.
- * Many (almost all) algorithms use one and improve itself by trying to minimize the loss function.
- * Some of them (maybe all) can be used as metrics to evaluate the model performance in the end (through ~~validation~~ data).

→ COST FUNCTION

- Loss function → For single training example
- Cost function → Average loss for entire dataset.

LOSS

CLASSIFICATION

- Binary Cross Entropy (Log-Loss)
- Hinge Loss
(For labels -1 & +1)

MULTI-CLASS

- Multi-class Cross Entropy
- Kullback-Liebler Divergence

REGRESSION

- Squared Error
- Absolute Error
- Huber Loss

CLASSIFICATION

Binary Cross Entropy (Log-Loss)

$$L = -y \cdot \log(p) - (1-y) \log(1-p) = \begin{cases} -\log(1-p) & \text{if } y=0 \\ -\log(p) & \text{if } y=1 \end{cases}$$

⇒ For calculating p:

$$\text{SIGMOID} = \frac{1}{1+e^{-z}} \quad \text{where } z: \text{a function of our input features}$$

Hinge loss

- * Primarily used with SVM with classes [-1, +1]

$$L = \max(0, 1 - y \cdot f(x)) \Rightarrow \text{for } y \in [-1, 1]$$

- * Not only penalizes wrong predictions but also the right predictions that are not confident.

- * Used when we want to make real-time decisions with not a laser-sharp focus on accuracy.

MULTICLASS CLASSIFICATION

Multiclass Cross-Entropy Function

→ A generalization of the Binary Cross-Entropy Loss.

$$\bullet L(x_i, y_i) = - \sum_{j=1}^C y_{ij} \cdot \log(p_{ij})$$

↳ where y_i one-hot encoded vector

$$\text{SOFTMAX} : \sigma(z)_i = \frac{e^{z_i}}{\sum_j e^{z_j}}$$

* is implemented through a neural network layer just before the output layer.

↳ The softmax layer must have the same number of nodes as the output layer.

Kullback - Liebler Divergence

→ A measure of how a probability distribution differs from another distribution. If zero → identical distributions

$$D_{KL}(P||Q) = \int P(x) \cdot \log\left(\frac{P(x)}{Q(x)}\right) \cdot dx \quad (\text{or } \sum \text{ for discrete})$$

↳ Expectation of logarithmic difference between P & Q in respect to P .

→ Not SYMMETRIC. So cannot be used as a distance metric.

① Minimize forwards KL $\rightarrow D_{KL}(P||Q) \rightarrow \text{Supervised}$

② Minimize backward KL $\rightarrow D_{KL}(Q||P) \rightarrow \text{Reinforcement}$

$$\rightarrow D_{KL}(P||Q) = - \sum_x (P(x) \cdot \log(Q(x)) - P(x) \cdot \log(P(x))) = H(P, Q) - H(P, P)$$

where $H(P, Q)$ is the cross-entropy of P & Q

$H(P, P)$ is the entropy of P .

REGRESSION

Squared Loss error

$J = (y - f(x))^2 \rightarrow$ cost function MSE

Absolute Error Loss

$$L = |y - p(x)| \rightarrow \text{cost function MAE}$$

Huber Loss

- Combines the best of MSE and MAE.
 - Quadratic for smaller errors and is linear otherwise.

$$L_f = \begin{cases} \frac{1}{2} (y - f(x))^2, & \text{if } |y - f(x)| \leq \delta \\ \delta |y - f(x)| - \frac{1}{2} \delta^2, & \text{otherwise} \end{cases}$$

⇒ Basically select δ and if Abs. Error is smaller than take Squared Error, else a function of Abs. Error & δ .

→ Used in Robust Regression / M-Estimation / Additive modelling

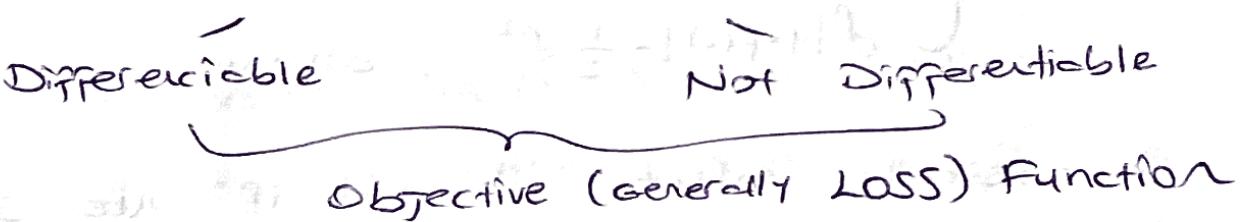
- tensorflow: keras, lasagne

OPTIMIZATION ALGORITHMS

- * Tries to find input parameters, weights etc. to a function that result in the minimum (maximum output) of the function.
- * With loss functions, ML algorithms use optimization algo's to minimize losses and find an optimal model.
 - ↳ Not specifically loss but COST function for the model. But COST function is an aggregation of the LOSS function, anyway.

MAIN BREAKDOWN

The function is continuous and can be differentiable at all points or not. This is an important factor.



DIFFERENTIABLE

1st ORDER : Slope of objective

↳ GRADIENT = Derivative of a multivariate continuous obj. func. (vector)

↳ PARTIAL DERIVATIVE = Element of the gradient vector

2nd ORDER : Rate at which derivative of the function changes

↳ HESSIAN MATRIX = Multivariate 2nd derivative

⇒ There are many ML algorithms (almost all that I know) use 1st order derivative for optimization. AS:

- Gradient Descent
- Momentum
- Adagrad
- RMSProp
- Adam

- Stochastic Gradient Descent
- Batch Gradient Descent
- Mini-batch Gradient Descent

Gradient Descent (core)

- We have a cost function differentiable and has local minima.
- Started with an initial value for model parameters, find the cost function, find the derivative (gradient) of the cost function.
- coefficient = coefficient - $(\alpha \cdot \Delta)$ where α = learning rate
- We go downhill as we know the directions by the gradient and we move coefficients in that direction as much as it is allowed by learning rate.

① Batch Gradient Descent - SLOW -

- As Gradient Descent is used for the entire dataset at each iteration → Batch Gradient

② Stochastic Gradient Descent - FAST -

- At each step, it takes a random record, uses it to update coefficients and repeats this process.
- Jumps a lot, overshooting → DECREASE LEARNING RATE

③ Mini-Batch Gradient Descent - MIDDLE -

- Takes a set (with smaller number than the train set) and take gradient for these data each step, update parameters and repeat this process.
- Choose mini-batch size wisely (generally 50-256)

NOTE : mini-batch size = 1 → Stochastic Gradient Descent

» Learning rate & Other parameter selection is important!

- * Gradient Descent is utilized in different forms in terms of selection of size in the batch.

⇒ However, there are lacks in these methods, and there are several optimizers built on top of these.

Momentum - with SGD

- * That helps accelerate SGD in the relevant direction and dampen oscillations.

- * Adds a fraction γ of the update vector of the past time step to the current vector:

$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta) \quad \gamma = \text{momentum (often 0.9)}$$

$$\theta = \theta - v_t$$

↳ With this, if gradient changes direction when compared to last round, velocity is slowed

↳ If it does not; velocity is increased

FOR RELEVANT PARAMETERS.

Nesterov Accelerated Gradient

By replacing θ with $\theta - \gamma v_{t-1}$:

$$v_t = \gamma v_{t-1} - \eta \nabla_{\theta} J(\theta - \gamma v_{t-1})$$

$$\theta = \theta - v_t$$

→ We give an idea of the future update, make the calculation with looking ahead to approximate future positions of our parameters instead of currents.

→ SMARTER dropping ball that has a notion of where it is going so that it knows to slow down before the hill slopes up again.

Adagrad

- It updates the learning rate for each parameter at each step.
 - * Low learning rates for \rightarrow frequently occurring features
 - * High learning rates for \rightarrow infrequent features
- Well suited for sparse data.

$$* g_{t,i} = \nabla_{\theta_i} J(\theta_{t,i}) \quad \text{at time } t, \text{ for parameter } i$$

$$* \theta_{t+1,i} = \theta_{t,i} - \eta g_{t,i} \quad \left\{ \text{SGD update} \right.$$

$$\hookrightarrow \theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{G_{t,i} + \epsilon}} g_{t,i}$$

where $G_{t,i} \sim$ past parameter gradients (sum of squares)
 $\epsilon \sim$ smoothing term to avoid divide by zero
 $(\sim 1e-8)$

- Accumulation of squared gradients in the denominator keeps growing and it leads to diminishing learning rate, which, in some cases, might result in model getting stuck.



Adadelta

- Built on Adagrad but with an attitude of reducing aggressive, monotonically decreasing learning rate.

- * Instead of accumulating all past squared gradients, Adadelta restricts the window of accumulated past gradients to some fixed size w .

→ more efficiency: Don't store w previous squared gradients. Store it decaying average of all past squared gradients.

$$\bullet E[g^2]_t = \gamma \cdot E[g^2]_{t-1} + (1-\gamma) g_t^2$$

$\gamma \sim$ momentumish

Then :

$$\Delta\theta_t = -\frac{\eta}{\sqrt{g_t^2 + \epsilon}} \odot g_t = \frac{-\eta}{\sqrt{E[g^2]_{t-1} + \epsilon}} \cdot g_t = \frac{-\eta}{\text{RMS}[g_t]} \cdot g_t$$

$$\rightarrow E[\Delta\theta^2]_t = \gamma E[\Delta\theta^2]_{t-1} + (1-\gamma) \Delta\theta_t^2$$

$$\hookrightarrow \text{RMS}[\Delta\theta]_t = \sqrt{E[\Delta\theta^2]_t + \epsilon}$$

* which we do not know! Need to be estimated

$$\Delta\theta_t = \frac{\text{RMS}[\Delta\theta]_{t-1}}{\text{RMS}[g_t]} \cdot g_t$$

$$\theta_{t+1} = \theta_t + \Delta\theta_t$$

→ Adadelta, we do not need to set a default learning rate as it has been eliminated from the update rule.

RMSprop

- Again, works on Adagrad's diminishing learning rate.
- RMSprop ~ IDENTICAL ~ to first update vector of ADADelta!

$$E[g^2]_t = 0.9 E[g^2]_{t-1} + 0.1 g_t^2$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} \cdot g_t$$

- $\gamma = 0.9$ by Hinton (developed RMSprop), while a good default value for learning rate ~ 0.001 .

Adam

- Adaptive learning rates for each parameter as well.
- exponentially decaying avg of past squared gradients
- + exponentially decaying avg. past gradients m_t , similar to momentum.

$$m_t = \beta_1 m_{t-1} + (1-\beta_1) g_t \quad \sim \text{mean gradients}$$

$$v_t = \beta_2 v_{t-1} + (1-\beta_2) g_t^2 \quad \sim \text{variance gradients}$$

⇒ Initialized as 0 vectors. Biased towards zero.

To Attack:

$$\hat{m}_t = \frac{m_t}{1-\beta_1^t}, \quad \hat{v}_t = \frac{v_t}{1-\beta_2^t}$$

$$\Rightarrow \left\{ \theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \cdot \hat{m}_t \right\}$$

$$\beta_1 = 0.9$$

$$\beta_2 = 0.999$$

$$\epsilon = 10^{-8}$$

} Proposed default values

Adamax

→ v_t factor was for l_2 -norm of past gradients

BUT it can be EXTENDED to l_p form.

$$u_t = \beta_2^\infty u_{t-1} + (1-\beta_2^\infty) |g_t|^\infty = \max(B_2 \cdot v_{t-1}, |g_t|)$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{u_t} \cdot \hat{m}_t$$

$$\eta = 0.002$$

$$\beta_1 = 0.9$$

$$\beta_2 = 0.999$$

} Proposed default values

N-Adam

- * Nestrov - Accelerated Adaptive Moment Estimation
- * Combines Adam & NAG.

$$g_t = \nabla_{\theta_t} J(\theta_t)$$

$$m_t = \gamma m_{t-1} + \eta g_t$$

$$\theta_{t+1} = \theta_t - m_t = \theta_t - (\gamma m_{t-1} + \eta g_t)$$

Also: $g_t = \nabla_{\theta_t} J(\theta_t - \gamma m_{t-1})$ From NAG.

↪ But don't apply twice! So use $J(\theta_t)$

Instead

$$m_t = \beta_1 m_{t-1} + (1-\beta_1) g_t$$

$$\hat{m}_t = \frac{m_t}{1-\beta_1^t}$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \cdot \hat{m}_t$$

↪ When replaced with its values:

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \left(\beta_1 \hat{m}_t + \frac{(1-\beta_1) g_t}{1-\beta_1^t} \right)$$

AMSGrad

→ Exponential averaging prevent learning rates to become infinitesimally small.

→ BUT, in settings where Adam converges to a suboptimal solution, some minibatches provide large & informative gradients, but they are RARE, exponential avs. diminishes their influence
 ↪ POOR CONVERGENCE.

- $v_t = \beta_2 v_{t-1} + (1-\beta_2) g_t^2$, $\hat{v}_t \Rightarrow NO!$

- $\hat{v}_t = \max(\hat{v}_{t-1}, v_t)$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \cdot m_t$$

SECOND ORDER OPTIMIZERS

→ Université Obj. Functions

- Newton's method
 - Secant method

→ Multivariate OLT- Functions

- ## • Quasi - Newton Methods

NON-DIFFERENTIAL OPTIMIZERS

→ Direct Algorithms

- Cyclic Coordinate
 - Powell's method
 - Hooke-Jeeves method

→ Stochastic Algorithms

- Simulated Annealing
 - Evolution Strategy
 - Cross-entropy method

→ Population Algorithms

- Genetic Algorithm
 - Differential Evolution
 - Particle Swarm Optimization

more advanced topics

CONCEPT OF ENTROPY

- * In Information Theory,
 - Entropy of a RV. is the average level of "info", "surprise" in its possible outcomes.

- * For X with outcomes x_i with prob. p_i

$$H(x) = - \sum_{i=1}^n p(x_i) \log(p(x_i)) \quad \text{ENTROPY}$$

- log base 2 bits

base e nats

base 10 dits

$$G(x) = 1 - \sum_{i=1}^n p^2(x_i) \quad \text{GINI}$$

- High Entropy : More uncertainty

- Low Entropy : More predictability

Information Gain

- * Concept of a decrease in entropy after splitting the data on a feature.

$$IG(D_{pf}) = I(D_p) - \sum_{j=1}^m \frac{N_j}{N} I(D_j)$$

- I = either entropy or gini
- D_j = dataset on the j^{th} split
- D_p = dataset on the parent node
- N = # of data in parent
- N_j = # of data in j^{th} child

⇒ Decision Trees split by looking at IG !

VALIDATION

- * We need to validate our trained model with a hold-out dataset to test how it performs outside-the-box.

Simple Train-Test Split

- * Split the data in two parts
 - ↳ use 1 to train
 - ↳ use other to test

* You can stratify!

- * PROBLEMATIC for MULTICLASS / MULTILABEL data
 - ↳ Stratified Shuffle Split

CROSS-VALIDATION

- ① To compare different models or different hyperparameters as a part of

- different model error stats & their deviations
- some model different hyperparameters' error stats & deviations

↳ Then, DECIDE which model is more powerful AND Stable.

- ② Small set of data and you cannot split a test data completely separate → LOSS OF INFORMATION

↳ Do CROSS-VALIDATION:

- * If the error is satisfying and

- * Error does not fluctuate a lot !

↳ Train a model on the whole dataset
(FINAL MODEL)

Leave - One - Out cross validation

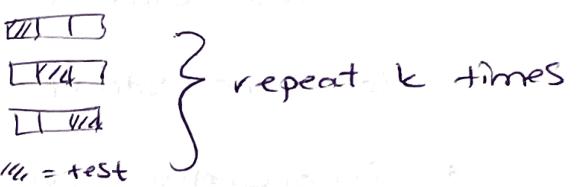
Training Data = $n-1$ }
 Test Data = 1 } n times repeat

↳ less bias

↳ computationally expensive

k-fold cross Validation

This time, data is folded into k parts, each k time:

- Train: k-1 parts }
 - Test: kth part }
- 
- repeat k times

⇒ When $k=n \rightarrow$ LOOCV!

* You can stratify to deal with imbalanced datasets

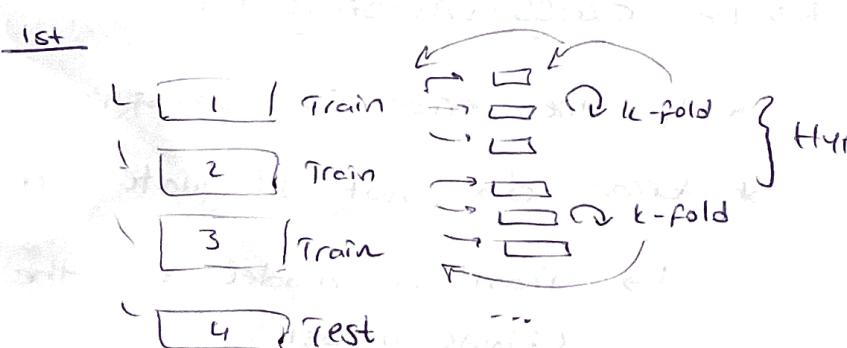
* You can also shuffle dataset each fold.

* WHAT IF:

* You have a small dataset where you have to build a model and tune your hyperparameters?

- NESTED k-fold CV

⇒ k-fold CV for hyperparameter optimization is NESTED INSIDE the k-fold CV for model selection.



Then you find best hyperparameter sets, REFIT them to all train 1st layer data

↳ Test those models → repeat

↳ REACH a optimized set with model stats.

→ But it might be costly in terms of the INCREASE in the NUMBER OF MODEL EVALUATIONS performed.

$$\hookrightarrow k \cdot n \cdot k$$

PYTHON IMPLEMENTATION

* sklearn.model_selection.

↳ train-test-split ($X, y, test_size = ?$,
shuffle = False / True,
random-state = ?,
stratify = array / None) = $x_{tr}, y_{tr}, x_{test}, y_{test}$

↳ Kfold ($n_splits = ?$,
shuffle = True / False,
random-state = ?) } Object

→ .split ($X, y = \text{None}, groups = \text{None}$) = train, test (index)

↳ StratifiedKFold ($n_splits = ?$,
shuffle = True / False,
random-state = ?) } Object

→ .split ($X, y, groups = \text{None}$) = train, test (index)

↳ cross-val-score (estimator,
 $X, y = \text{None}, groups = \text{None}$,
groups = None, cv = None, n-jobs = None / -1,
scoring = 'str' or score callable (est., X, y),
cv = None (s),
n-jobs = None / -1, verbose = 0,
fit_params = None,
error_score = 'raise', np.nan) → Scores

↳ RepeatedKFold ($n_splits, n_repeats, random_state$)

↳ LeaveOneOut()

↳ GroupedKFold(...)

ALSO

about n-grams and how to implement them.

* `Sklearn.cross-validation`: ~~cross-validation~~ is a module that provides various cross-validation strategies for evaluation.

L Looks like some functions are defined here as well.

⇒ Write code for further xd

BIAS - VARIANCE TRADEOFF

BIAS: The error that is introduced by modelling a real-life problem.

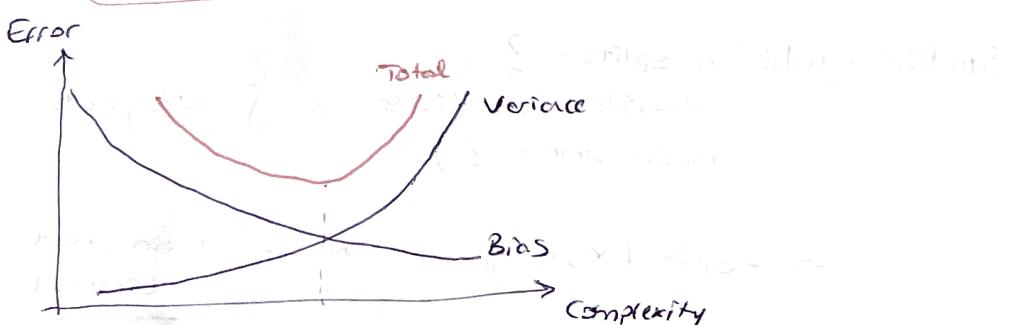
→ Our prediction vs truth difference

→ High bias = Simple model (over!!!)

VARIANCE: Variability of model prediction for given data.

→ High variance = Complex model that overfits

$$\text{Err}(x) = \text{Bias}^2 + \text{Variance} + \epsilon \quad \rightarrow \text{minimize!}$$



* Low bias good fit but as it goes too low as the model fits more and more, variance increases.

↪ Find the optimal by testing your model with hold-out-sets and compare metrics.

↪ Fix your model / hyperparameters if NEEDED!

⇒ FIND THE BEST BALANCE!

⇒ NO OVERFIT / NO UNDERFIT

SUPERVISED MODELS

- * Data is labeled and labels are utilized during the model parameter optimization process.
 - ↳ CLASSIFICATION: Target is finite and unordered set.
 - ↳ REGRESSION: Target is continuous and quantitative.

CLASSIFICATION

- ① Binary Classification
 - ② Multi-class Classification
 - ③ Multi-label Classification
- } x Also combination of these

LOGISTIC REGRESSION

→ It does work as binary classification model, but also gives probabilities to be used as crasing.

→ Hypothesis $z = w^T x + b$

$$h(x) = \text{sigmoid}(z), \text{ where sigmoid}(z) = \frac{1}{1+e^{-z}}$$

$$= P(Y=1 | X; \theta)$$

→ Data is fit into linear regression model, which then be acted upon by a logistic function predicting the target categorical dependent variable



→ S-shaped curve fits data that gives $(0,1)$ probability function.

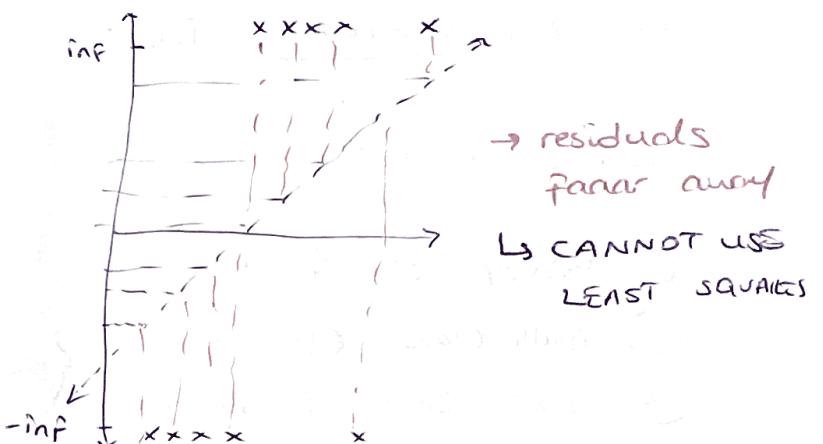
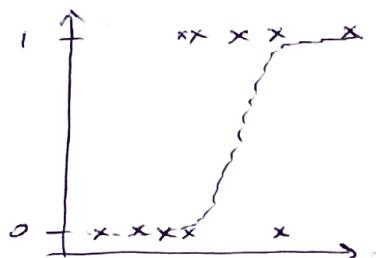
→ Uses maximum likelihood and calculates it at each iteration until maximizing it.

→ Logit = $\log\left(\frac{P}{1-P}\right)$, function reverses it into linear line prob.

→ Intercepts are in linear function. | (after logit)
Coefficients

→ When categorical, log(odds) used again
also difference in log(odds) s.

* COEFFS are in terms of log(odds)



- Instead we project these points onto the line and find log(odds) correspondents.
- Transform coordinate by sigmoid function
- Likelihoods are predicted probabilities (P) on the line.

⇒ Take the products of likelihoods for each class:

- Positive class $\pi(P)$
- Negative class $\pi(1-P)$

BUT: Log-Likelihoods

$$\ln \log(P)$$

→ Rotate the line → go to P_s → max log-likelihood

* NO OUTLIERS

* NO MULTICOLLINEARITY IN PREDICTORS

→ Uses log-loss function so that gradient descent can be used to find a minimum in convex space.

PYTHON

* Sklearn. linear-model. Logistic Regression

```
(penalty = 'l1', 'l2', 'elasticnet', 'None',  
dual = False,  
C = 1.0,  
class_weight = None,  
random_state = ?,  
solver = newton-cg / libfgs / liblinear / sag / saga,  
multi-class = 'auto' )
```

- coef_
- intercept_
- n-Features-in-
- Feature-names-in-
- n-iter-
- fit(x, y)
- predict(x)
- predict_proba(x)

Naive Bayes Classifier

Based on Bayes' Theorem, with the assumption of independence between every pair of Features.

- Extremely Fast.
- Can work with small set of data.
- Document classification and spam filtering
- * Beginning of this NOTEBOOK!

GAUSSIAN

A priori probability of class A ($P(A)$) is used. Also $B \rightarrow P(B)$

$$P(A) \times L(X_1=x_1 | A) \times L(X_2=x_2 | A) \dots = P(A | X_1=x_1, \dots)$$

$$P(B) \times L(X_1=x_1 | B) \times L(X_2=x_2 | B) \dots = P(B | X_1=x_1, \dots)$$

where likelihoods are pdf values of Normal Gaussian of features.

→ MAY USE LOG() function $\log() + \log(1 - \dots)$

* WHICH SCORE higher → classified

PYTHON

* Sklearn, naive-bayes. GaussianNB

()

- fit(x, y)
- predict(x)

↳ can be easily writable by self.

K - Nearest Neighbors

→ Decide a number k which will be the nearest neighbors number to ask their vote.
The MAJORITY class voted is the data's class.

- * Simple, used for nonlinear data
- * Gives high accuracy but there are better models
- * It stores ALL training data and calculates distances then orders and decides.

* Decide k with cross validation

Python

* sklearn.neighbors.KNeighborsClassifier

```
(algorithm='auto',
leaf_size=30,
metric='minkowski',
metric_params=None,
n_neighbors=5,
p=2,
weights='uniform')
```

- fit(x, y)
- predict(x)

DECISION TREE

- Creates rule patterns.
- Top-down, greedy, recursive binary splitting.
- * It uses previous concepts of gini and entropy to calculate impurities, ergo information gain on splits.
 - ↳ Maximizing IG is the selection method for feature and cut points.
- No need to normalize data.
- can be used as preprocessing tool for missing value prediction.
- No assumptions about distribution.
- Might overfit on noisy data OR high dimensional scarce data
- Biased with imbalance dataset

PRUNING

- We can grow the tree as the improvement criterion exceeds a threshold.
- * BUT SHORT SIGHTED: an unimportant split can be followed by an important one.
- First, build the model.
- Then prune with COST COMPLEXITY PRUNING

$$\sum_{m=1}^{|T|} \sum_{i:x_i} (y_i - \hat{y}_{R_m})^2 + \alpha |T| L$$

→ Eliminate weakest link

PYTHON

- * `sklearn.tree.DecisionTreeClassifier(`
- `criterion = 'gini' / 'entropy'`,
- `splitter = 'best'`
- `max_depth = None`,
- `min_samples_split = 2 (number OR fraction)`)

}
 - feature-importances -
 - max-features -
 - n-features -
 - feature-names-in
 ...

- fit(X, y)
- predict(X)

* `sklearn.tree.DecisionTreeClassifier()`

- cost-complexity-pruning-path (X, y)

- ccp-alphas } plot
- impurities

↳ To select your classifier, compare the above. OR
 ↳ Build a classifier set with $\text{ccp_alpha} = \text{ccp_alpha}_i$
 ↳ Run and calculate accuracy for training and test sets. OR CV!

Regressor Tree

* Instead of majority count, takes the mean of the leaf to predict.

* `DecisionTreeRegressor()`
 $\text{criterion} = \text{squared-error / mse / absolute-error/mae}$, ...)

But ONE TREE IS NOT THE BEST!

BAGGING ~ Bootstrap Aggregating ~

* Bootstrapping = Resampling with replacement

→ Here, same as the overall sample size, create B bootstrapped B samples.

→ Fit each one a tree.

→ Aggregate each tree outcome by majority vote among trees.

* Reduces variance ...

→ Each weak learner is constructed parallelly, not one top of each other.

→ Standard = Bootstrapped Sample Size = Sample size

Out-of-Bag (OOB) error

↳ If you have 1000 data points, 200 will be in the bag.
Since sampling is done by replacement, some records occur more than once, some are never in the bag.

~ Appx. 2/3 - 1/3 is Out-of-Bag.

→ Train the model on the bag and predict the records that are OOB and measure the metric of interest.

↳ The average of all metrics for each bag is the model metric. We also have deviation

≈ kind of ~~CROSS-VALIDATION~~ ~~cross-validation~~

Python

* sklearn.ensemble.BaggingClassifier

base_estimator = DecisionTreeClassifier,

n_estimators = # bases,

max_samples = bag size (can be %),

max_features = feature set size (can be %),

bootstrap = True / False,

bootstrap_features = True / False,

oob_score = True / False,

random_state = ?) . . .

base-estimator-
n-features-
n-features-in
feature-classes-in
estimators-
estimators-Samples-
estimators-Features-
classes-
oob-scores-
oob-decision-function-

- fit(x, y)

- predict(x)

- predict_proba(x)

⇒ You can create different base models etc...

RANDOM FOREST

Random Forest (RF)

- Built on the idea of bagging.
- Takes many trees from bootstrapped samples but this time, at each sample, a random sample of m predictor is chosen as candidates from full set of P .
- conventionally $m \sim \sqrt{P}$
- * If there are strong predictors among some moderate ones (in real life always :)), bagging will select strong ones as the first split at each learner.
 - ↳ use a random subset of features.

Implementation Using DECORRELATED TREES

PYTHON

`sklearn.ensemble.RandomForestClassifier()`

```

n_estimators = 100,
criterion = gini / entropy,
max_depth = None,
min_samples_split = 2,
min_samples_leaf = 1,
max_features = auto,
max_leaf_nodes = None,
bootstrap = True,
oob_score = False,
class_weight = 'balanced',
ccp_alpha = 0.0,
max_samples = None
  
```

In addition to bagging:
 · feature-importances-

- `fit(x, y)`
- `predict(x)`
- `predict_proba(x)`

ISOLATION FOREST

* Built on the basis of decision trees, mostly for outlier detection.

→ If you can separate a sample in the early cuts then it means it is more away from other values, potential OUTLIERS! If a sample is in a deep leaf of a tree, it means it was hard to distinguish it from others (regular values)

$$\text{Anomaly Score} = S(x, n) = 2 - \frac{E(h(x))}{c(n)}$$

$h(x)$: path length of obs. x

$c(n)$: avg path length of unsuccessful search

Search in a Binary Search tree

n = number of external nodes

→ After fitting the algorithm, it calculates anomaly scores for all data points and with thresholds it plots them.

Python

* Sklearn.ensemble.IsolationForest()

$n_estimators = 100$, } previous tree based returns
 $\max_samples = \text{auto}$, } +
 $\max_contamination = \text{auto}$, } offset -
 $\max_features = 1.0$, }
 $\text{bootstrap} = \text{false}$, }
 $\text{random_state} = \text{None}$) }
 fit(X) → $+/-$
 fit-predict(X) → $+1/-1$
 predict(X) → $+1/-1$
 decision-function(X) → anomaly-scores

⇒ It is somehow unsupervised!

↳ No labels as outlier/no outlier...

EXTENDED ISOLATION FOREST

→ Again UNSUPERVISED, same mechanism but extended.

→ Isolation Forest has only vertical/horizontal cuts.

Therefore if there are multi-modes in the data then may mislead.

↳ NEVER USE SIMPLE ISOLATION unless the problem is too simple and obvious.

↳ ALTHOUGH, if you are close to a simple feature each time (1-dimension) then OKAY.

PROCEDURE

- * Random slope for the branch cut
- * Random intercept chosen from the range of available values from the training data

SELECTED

* Makes tailored hyperplanes & better results

Python

```
eif.eForest(x, n_estimators=n_trees,  
           sample_size=sample_size,  
           ExtensionLevel=extension_level)
```

By ExtensionLevel we can create Extended Iso.

BOOSTING

* Trees are grown sequentially, each tree grown using information from previously grown trees.

→ DOES NOT HAVE TO BE SET OF TREES. H can take other models as well.

PROCEDURE

1 → Start with initial equal weights to all samples.

2 → Take randomly one, fit the tree / model, store predictions & RESIDUALS

3 → Give more weight to samples with higher residuals

4 → Repeat from 2.

$$\text{Amount of say} = \frac{1}{2} \log \left(\frac{1 - \text{Total Error}}{\text{Total Error}} \right) \rightarrow \text{sample weights}$$

ADAPTIVE BOOSTING

* In ADABOOST, trees are only stumps contrary to RF where tree lengths must be long and varied among weak classifiers.

o STUMPS = 1 node binary splits

PROCEDURE (numbers example)

1. Divide data to length of data to give initial sample weights.

2. Make the first stump by finding the variable that does the best job. For each variable - split - find gini / entropy & SELECT the best one.

3. Sum of the weights classified as wrong → TOTAL ERROR = $\frac{1}{8}$

$$4. \text{ Calculate Amount of say} = \frac{1}{2} \log \left(\frac{1 - \frac{1}{8}}{\frac{1}{8}} \right) = \frac{1}{2} \log(7) = 0.97$$

$$5. \text{ New Sample weight} = \text{Sample weight} \times e^{\text{amount of say}} = \frac{1}{8} \cdot e^{0.97} = 0.33 \quad (\text{INCORRECT})$$

$$6. \text{ New Sample weight} = \text{Sample weight} \times e^{-\text{amount of say}} = \frac{1}{8} \cdot e^{-0.97} = 0.05 \quad (\text{CORRECT})$$

7. Normalize sample weights as summing up to 1.

8. Pick a random number and use sample weights as cdf for sample selection

9. Bootstrap sample as the same size with original with this random selection process.

Notice weight-increased samples of incorrectly predicted will appear more in bootstrapped Samples.

10. Repeat this process until n-estimators are reached.

⇒ Amount of says calculated are also used as vote weights of the stumps in the final decision.

⇒ class P votes \sum amount of says
class N votes \sum amount of says } compare ⇒ CLASS

PYTHON

* sklearn.ensemble.AdaBoostClassifier(

base-estimator = DecisionTreeClassifier(max-depth=1),

n-estimators = 50,

learning-rate = 1.0,

algorithm = SAMME.R / SAMME,

random-state = None)

Addition:
.estimator-weights
.estimator-errors
.feature-importances

- fit(x,y)

- predict(x)

- predict-proba(x)

GRADIENT BOOSTING

→ The idea is the same. But this time, weights are not the manipulation tool.

→ A LOSS function is defined as it will give us the residuals for our predictions.

→ Then, the next weak learner tries to fit the model to those residuals.

PROCEDURE

1. Define a differentiable LOSS FUNCTION. that IS comparable with base estimator.
2. Initialize your model $F(x) = \arg\min_{\gamma} \sum_{i=1}^n L(y_i, \gamma)$
3. For $m=0$ to M :
 - Compute pseudo-residuals $r_{im} = - \left[\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]$
 - Fit a base learner $h_m(x)$ to residuals using the train set $\{(x_i, r_{im})\}_{i=1}^n$
 - Compute multiplier $\gamma_m = \arg\min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i))$
 - Update the model $F_m(x) = F_{m-1}(x) + \gamma_m h_m(x)$
4. Repeat 3 until n-estimators reached.

⇒ By this, we adjust our rate in residual-fitted model whilst minimizing the loss function; AND we add it to the previous model, basically ADJUSTING the previous model more accordingly to its MISTAKES.

⇒ Does not have to be Stumps! But different trees

Combined.

⇒ Named as Gradient since it uses minimizing loss function by gradient descent.

Python

* `sklearn.ensemble.GradientBoostingClassifier()`

loss = deviance / exponential,

learning-rate = 0.1 ,

n-estimators = 100 ,

subsample = 1.0 (smaller than 1.0 → SED)

criterion = friedman_mse / squared_error, mse, mae,

min_samples_split = 2 ,

min_samples_leaf = 1 ,

```

min_weight_fraction_leaf = 0.0,
max_depth = 3,
min_impurity_decrease = 0.0,
random_state = None,
max_features = None / auto / sqrt / log2,
verbose = 0 (prints each tree),
max_leaf_nodes = None,
validation_fraction = 0.1,
n_iter_no_change = None, // early stop //
tol = 1e-4,
ccp_alpha = 0.0
)

```

} oob_improvement -
train_score -
loss -
... other scores .

XGBoost

- Gradient Boosting is SLOW since it is SEQUENTIAL.
- XGBOOST is the same but further developed to solve performance / efficiency issues.

- * parallelization of tree construction using all CPUs
- * Distributed computing using cluster of machines
- * Out of Core computing that data not fitting into memory
- * Cache Optimization to make the best use of hardware

xgboost.XGBClassifier(

```

n_estimators
use_label_encoder
max_dept
learning_rate
verbosity
objective = 'binary:logistic',
booster
tree_method = 'auto'
gamma = minimum loss reduction required,
)

```

+ lambda = L2 reg. (= 1)
alpha = L1 reg. (= 0)

```

.fit(x, y)
.predict(x)
.predict_proba(x)

```

Objective

- o binary: logistic \rightarrow log regression for binary prediction \rightarrow prob.
- o multi = softmax \rightarrow multiclass classification using softmax \rightarrow class \hookrightarrow num-class
- o multi = softprob \rightarrow Some ass softmax but returns probs.

Errors (eval-metric)

- rmse
- mae
- logloss (neg. log-likelihood) \rightarrow multiclass
- error
- merror (multiclass)
- mlogloss (multiclass)
- auc

LIGHT GBM

- Microsoft released in 2017
- Others split depth-wise whereas Light GBM splits leaf-wise

* Finding the best split for each leaf is really hard. Naively done $\rightarrow O(n.p)$ complexity.

TECHNOLOGIES

* Gradient Based One Side Sampling

- Sort the instances acc. to gradients in descending order
- Select the top $\alpha \times 100\%$ instances
- Randomly sample $b \times 100\%$ from the rest of the data
 \Rightarrow Putting more weight on untrained without disrupting data distribution

* Exclusive Feature Bundling

- Construct measure of conflict between features (based on overlapping non-zero values)
- Sort the features by count of non-zero instances in descending
- Loop over the ordered list of features and assign the feature to an existing bundle ($\text{conflict} < \text{threshold}$)
 new bundle ($\text{conflict} > \text{threshold}$)

→ This helps in dimension reduction (somehow) and speed up the process.

→ Some features are never non-zero together.

Python

lightgbm.LGBMClassifier(

boosting-type = 'gbdt' ,

num-leaves = 31

max-depth = -1

learning-rate = 0.1

n-estimators = 100

sumsample-for-bin = 200,000,

objective = None (default),

reg-alpha = L1 ,

reg-lambda = L2 ,

)

- objective -
- evals-result -
- booster -
- best-iteration -
- best-score -

.fit(x, y)

.predict(x)

.predict-proba(x)

HYPERTPARAMETER TUNING

RULES

- * Know which model you are working with and what that might suffer from.
↳ Depth, criterion, n-features etc.
- * Know your data and its size, have good propositions for hyperparameters
- * Create a dictionary (set) of variables & their proposed values / by variables = hyperparameters /
↳ Run a CV / OOB on data whilst testing different sets of hyperparameters.
- ↳ Know your SEARCH Algorithms

GRID SEARCH

- You can construct a dict with list of proposed values as values and hyperparameters as the keys
- * Then construct your grid search with CV by python with loops etc.
- Already constructed functions = GRIDSEARCHCV()
- * sklearn.model_selection.GridSearchCV(

```

estimator = ...,
param_grid = dict(...),
scoring = 'compatible score name',
cv = 5,
verbose = 1,
n_jobs = -1)
    
```

- fit(x, y)
- get_params()
- inverse_transform(x)
- predict(x)
- predict_proba(x)

{ best }

- cv - results -
- best - estimator -
- best - score -
- best - params -
- best - index -
- scorer -
- n - splits -
- multimetric -

RANDOM SEARCH

- * Create a grid and select random combinations among these:
- * `sklearn.model_selection.RandomizedSearchCV(`

```
    estimator = None,                                # Estimator
    param_distributions = dict / list of dicts,     # Distributions
    n_iter = 10,                                     # Iterations
    scoring = None,                                   # Scoring metric
    n_jobs = None,                                    # Number of parallel jobs
    refit = True,                                     # Refit model
    cv = None(s),                                     # Cross-validation
    verbose = int,                                    # Verbosity level
    random_state = None)
```

Some Grid.

→ Can take whole distributions as parameter value ranges in the dict.

BAYESIAN OPTIMIZATION

- * `Hyperopt.fmin()`

→ It uses bayesian probability to find the minimum of a function (loss)

```
(V) objective_function = loss
      domain_space = range of inputs, grid
      optimization_algorithm = "fmin"
```

→ Better than others.

SUPPORT VECTOR MACHINES

- Can be applied to a variety of settings inc. binary classification, multi-nominal classification, regression.
- Separating hyperplanes

PROCEDURE

1. Fit a $n-1$ hyperplane to an n -dimension space
2. Measure perpendicular distance to the hyperplane of each data
3. Take the smallest distance \rightarrow margin
4. Find a hyperplane that maximizes the margin.

↳ works if separable classes

What if not-separable classes

SOFT MARGINS

- * It allows some points being on the ~~right~~ side of the hyperplane.
- * Penalty term C:
 - Larger, more sophisticated models.
 - If kernel not changed, then a wider margin comes with larger C (penalty for misclassification) $\rightarrow \uparrow$ support vectors

$$\text{Max } M \quad \text{s.t. } \sum_1^P \beta_j^2 = 1$$

$$y_i(\beta_0 + \beta_1 x_{i1} + \dots + \beta_P x_{ip}) > M(1 - \epsilon_i)$$

$$\epsilon_i \geq 0, \sum_1^n \epsilon_i \leq C,$$

KERNEL TRICK

- * For finding non-linear decision boundaries.
- * change kernel \rightarrow
 - linear
 - poly
 - rbf
 - sigmoid
 - precomputed

RBF

$$\phi(x, \text{center}) = \exp(-\gamma \|x - \text{center}\|^2)$$

→ Transformer by measuring the distance between all other dots to a specific/dots-centers.

γ = rate of influence of new features
 ↳ more, more wissly

KERNEL : Data is projected/transformed into more ordered, standardized forms. During this we do not store new coordinates. Except, we store:

- Distances :

$$\langle x_i, x_j \rangle = \sum_1^p x_{ij} x_{i,j} \quad (\text{inner product})$$

$$f(x) = \beta_0 + \sum_{i=1}^n \alpha_i \langle x, x_i \rangle$$

↳ It turns out many α_i is 0.

↳ S = Support set of i, $\alpha_i > 0$.

Linear $\rightarrow x_i \circ x_j$

Polynomial $\rightarrow (1 + x_i \circ x_j)^d$

Radial $\rightarrow \exp(-\gamma \|x_i - x_j\|^2)$

Sigmoid $\rightarrow \tanh(\alpha x^T y + c)$

PYTHON

* sklearn.svm.SVC (

$c = 1.0,$

$\text{kernel} = 'rbf'$,

$\text{degree} = \text{if poly},$

$\text{gamma} = \text{'scale'}$,

$\text{coef0} = \text{if poly/sigmoid},$

$\text{shrinking} = \text{True},$

$\text{probability} = \text{False},$

$\text{tol} = 1e-3,$

$\text{--- decision-function-shape= ovr/ovo})$

• Support -
 • support-vectors -

• n-support -

• probA - / .probB -

• shape-fit -
 - ---

• fit(X, Y)

• predict(x)

SI UNSUPERVISED MACHINE LEARNING

↳ Goal: Discovering useful patterns in data

- * No labels / finding patterns intrinsic in the data

→ Parametric: Gaussian Mixture Models / Exp.-Max. Algos

→ Non-parametric: As we know it, dist-free

- * CLUSTERING

↳ Exclusive, Overlapping, Hierarchical, Probabilistic
(K-Means) (Fuzzy-k-Means) (Hierarchical) (Mixture of Gaussians)

K-MEANS CLUSTERING

→ This algorithm basically finds k cluster centroids and assign points to the cluster with the minimum distance.

• PROCEDURE • (Greedy)

① Randomly select k centroids

② Associate each point to the closest centroid

③ Minimize $J = \sum_{j=1}^k \sum_{i=1}^n \|x_i^{(j)} - c_j\|^2$

↳ distance btwn a data point x_i and the cluster centre c_j

④ Re-calculate cluster centers:

$$v_i = \left(\frac{1}{c_i}\right) \cdot \sum_{j=1}^{c_i} x_i \quad \text{where } c_i = \# \text{ of points in cluster } i$$

⑤ Repeat until no centroid is moved.

• - Initialization is not specified.

• - Set of samples closest to v_i is empty, so it cannot be updated. → Handle at implementation

→ use Elbow-method to decide K.

→ Scaling might be useful in improving K-means.

⚠️ → You cannot have CATEGORICAL variables

They should be NUMERIC!

$\sim O(\# \text{iters} * \# \text{clusters} * \# \text{instances} * \# \text{dimensions})$

! SCALE !

(related to **PYTHON**)

• `scipy.cluster.vq.kmeans()`

obs = data,

k-or-guess = ~~None~~, init = ~~random~~ ~~int~~,

miniter = 20, max_iter = ~~100~~ ~~max~~ ~~100~~,

thresh = $1e-05$,

check_finite = True,

seed = ~~None~~)

= centroids, labels

* df['labels'] = vq(data, centroids)

* `sklearn` → whiten(`data`) → variation reduction

• `scipy.cluster.KMeans()` { cluster-centers - }

{ n-clusters = 8, n-init = 10, max_iter = 300, } { labels - } { inertia - } { n-iter } { n-features } { feature-names-in }

* `MiniBatchKMeans()` ✓

• `fit(x)`

• `fit_predict(x)`

• `predict(x)`

• `score(x)`

HIERARCHICAL CLUSTERING

This is a method that starts with considering each data point as a cluster and then goes until there is only one cluster.

PROCEDURE

- ① N clusters start.
- ② Find the most closest clusters & merge them into one. \rightarrow + Less clusters \leftarrow
- ③ Compute distances and REPEAT 2 until done.

DENDROGRAM

- Visualization for hierarchical clustering.
- Best to decide on number of clusters
 - * HORIZONTAL LINE ON DENDROGRAM.

PYTHON

- Scipy.cluster.hierarchy

\hookrightarrow linkage(

```
y = data,  
method = 'single',  
metric = 'euclidean',  
optimal_ordering = False) = Z  
=====
```

\hookrightarrow dendrogram(Z) \rightarrow image
.plt.show()

\hookrightarrow fclusters(Z,

```
criterion = 'maxclust'  
numclusters = )
```

DIMENSIONALITY REDUCTION

→ When we have too many features (even small number can be too many), we may have to reduce dimension for avoiding information loss (clustering) or overfitting.

FEATURE ELIMINATION

- Variance Inflation Factor = Elimination method.
- Univariate Feature Selection = One-by-one ~ SFA
- Recursive Feature Elimination = fits a model, eliminate k weakest features.
- Model Based Feature Selection = feature importances etc.

* Lossing some information since features are directly dropped without any INFO EXTRACTION!

PRINCIPAL COMPONENT ANALYSIS

- Feature Extraction & Dimensionality Reduction Technique
- Might be used as a tool that takes features and get them into 2/3D visualization.

PROCEDURE

- ① Calculate a matrix that summarizes how our variables are all relate to one another, with direction and magnitude
- ② Transform our original data to align with these important directions.
 - ↳ Do it for all features.
 - ↳ Orthogonal directions.

- X matrix, each column is standardized to have 0 mean. (-mean for sure, but standardizes depends on your selection)
- Calculate covariance matrix = $Z^T Z$
- Calculate the eigenvectors by $Z^T Z \rightarrow PDP^{-1} \rightarrow \underline{P}$
- Take eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_p$, sort them in descending: $\rightarrow \underline{P^*}$ (sorted)
- Calculate $Z^* = ZP^*$ is a centered/standardized version of X . but now each obs. is a combination of the original variables, weights = determined by eigenvector

↳ EACH COLUMN OF Z^* INDEPENDENT!

↳ LINEAR INDEPENDENCE OF VARIABLES!

HOW MANY PCs?

- Arbitrarily Select. (\sim visualization, 2D, 3D)
- Proportion of Variance Explained
 - Threshold, add features until it is reached.
- Scree Plot
 - ↳ Variance \rightarrow elbow
 - ↳ Cumulative \rightarrow { again prop. of var. explained }

* PCR (Principal Component Regression)

* Kernel PCA

PYTHON

* `sklearn.decomposition.PCA()`

n_components = None,

copy = True,

whiten = False,

svd_solver = 'auto',)

• components -

• explained - variance -

• explained - variance - ratio -

• singular - values -

• mean -

• n - components -

• n - features -

• n - samples -

• fit(x) → finds - weighted -

• fit_transform(x)

• transform(x)

• inverse_transform(x) → original

+ Kernel PCA()

Space PCA()

Truncated SVD()

IncrementalPCA()

④

t-SNE

* t-distributed Stochastic Neighbor Embedding

→ PCA may require data in a linear relationship. If we have non-linear relationships → t-SNE.

→ It calculates a similarity measure based on the distance between points instead of trying to maximize variance.

PROCEDURE

① Measure similarities btw points in the high dimension.

For each data point (x_i) we center a Gaussian dist. over that point. Then we measure the density of all points (x_j) under that Gaussian.

↳ Renormalize for all points.

↳ Gives P_{ij} for all points ~ similarities

- use perplexity to adjust variance.

- ② Use Standardized (cauchy) dist. with n-1 d.o.f. This gives us:
↳ Second set of probs (Q_{ij}) in low dimensional space.

- ③ We want Q_{ij} to reflect P_{ij} as best as possible.

↳ Kullback-Liebler (KL) Divergence

- * Use Gradient Descent to minimize KL cost function

! Cannot extend to new samples.

! Each run different results but relative position is same.

PYTHON

- * `sklearn.manifold.TSNE()`

n-components = 2,
Perplexity = 30.0,
early-exaggeration = 12.0,
learning-rate = 200.0,
n-iter = 1000,
metric = 'euclidean',
init = 'random',
...)

} embedding -
• kl-divergence -
• n-features-in -
• feature-names-in
• n-iter -

• fit(x)
• fit-transform(x) ! NO: .transform() X

- * `sklearn.metrics.silhouette_score`