

IBM MACHINE LEARNING

Supervised Regression

$$y_p = f(\beta, x)$$

observations - row
features - column

$x \equiv \text{inputs}$

$y_p \equiv \text{predictions}$

$\beta \equiv \text{parameters of model or weights}$

hyperparameter - defines models, set during training

$J(y, y_p) \equiv \text{loss}$

Update: use x and y , choose β to minimize loss J

Interpretate as Predictive approach.

Linear Regression:

$$y_p(x) = \beta_0 + \beta_1 x$$

residuals, error $\sum ((\beta_0 + \beta_1 x_{\text{obs}}^{(i)}) - y_{\text{obs}}^{(i)})^2 \rightarrow \text{L2 norm}$

$$\text{minimize error, } J(\beta_0, \beta_1) = \frac{1}{2m} \sum ((y_{\text{pred}}^{(i)}) - (y_{\text{obs}}^{(i)}))^2$$

$$SSE = \sum (y_p^{(i)}) - y_{\text{obs}}^{(i)})^2$$

$$TSS = \sum (\bar{y}_{\text{obs}} - y_{\text{obs}}^{(i)})^2$$

$$R^2 = 1 - \frac{SSE}{TSS}$$

$$\text{boxcox}(y_i) = \frac{y_i^2 - 1}{\lambda}$$

scipy.stats.boxcox
scipy.stats.mstats.normtest

pf = PolynomialFeatures (degree= n, include_bias=False)
X_pf = pf.fit_transform(X)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state)

ss = StandardScaler()
X_train_ss = ss.fit_transform(X_train)
X_test_ss = ss.transform(X_test)

lr = LinearRegression()
lr.fit(X_train_ss, y_train)
y_pred = lr.predict(X_test_ss)
r2_score(y_test, y_pred)

Training $\xrightarrow[X_train]{y_train}$ model(X_train, y_train).fit() \rightarrow model

Test $\xrightarrow[X_test]{}$ model.predict(X_test) \rightarrow y_pred
 $\xrightarrow{}$ error-metric(y_test, y_pred) \rightarrow error/metric

train_test_split()
ShuffleSplit()
StratifiedShuffleSplit()

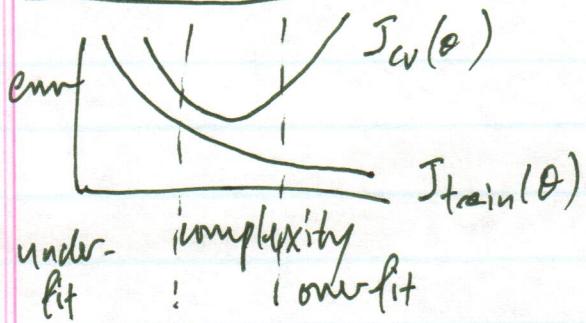
OneHotEncoder()
pd.get_dummies()

```

ax=sns.scatterplot(y-test, y-pred)
ax.set(xlabel=' ', ylabel=' ', title=' ')

```

Cross Validation:



K-fold cross validate
Leave one out CV
Stratified CV

```

cvn_val = cross_val_score(model, X_data, y_data, cv=k,
                           scoring='neg_mean_squared_error')

```

```
kf = KFold(shuffle=True, random_state=None, n_splits=n)
```

for train_index, test_index in kf.split(X):

```

X-train, X-test, y-train, y-test = (X.iloc[train_index, :],
                                      X.iloc[test_index, :],
                                      y[train_index],
                                      y[test_index])

```

lr.fit(X-train, y-train)

y-pred = lr.predict(X-test)

score = r2_score(y-test.values, y-pred)

scores.append(score)

```

estimator = Pipeline([('scaler', StandardScaler()),
                      ('regressor', LinearRegression())])

```

estimator.fit(X-train, y-train),

estimator.predict(X-test)

prediction = cross_val_predict(estimator, X, y, cv=kf)

r2_score(y, prediction)

~~np.mean(prediction)~~

Lasso (alpha = α)

smaller α , more complex, less regularized

• get_feature_names (= X.columns)

estimator.named_steps["regression"].coef_

Ridge (alpha = α)

smaller α , more complex, less regularized

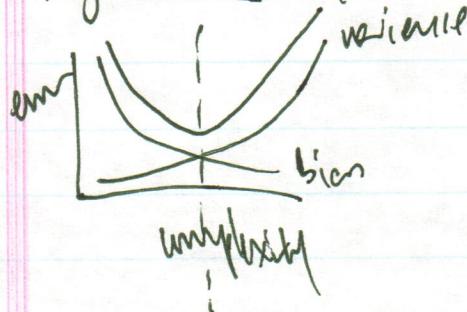
param = {
 'pp_degrees': [1, 2, 3],
 'ridge_alpha': np.geomspace(4, 20, 30)
}

grid = GridSearchCV(estimator, param, cv=kf)

grid.fit(X, y), grid.predict(X)

grid.best_score_, grid.best_params_, grid.best_estimator_.named_steps,
grid.cv_results_, .coef_

Regularization: total error



$$M(w) + \lambda R(w)$$

$M(w)$: model error

λ : regularize parameter / penalty

$R(w)$: regularize function

more regularization, larger λ , less complex model, more bias
less regularization, smaller λ , more complex model, more variance

if overfit, too high variance, regularize reduce error, variance.

Ridge (L2):

$$J(\beta_0, \beta_j) = \sum_i (y_i - \beta_0 - \sum_j \beta_j x_{ij})^2 + \lambda \sum_j \beta_j^2 = \text{residual sum of squares} + \lambda \sum_j \beta_j^2$$

- λ applied to β_j
- shrinking effect of β towards 0
- improves bias, reduces variance

$$J(\beta_0, \beta_j) = \frac{1}{2m} \sum_i ((\beta_0 + \beta_j x_{obs}^{(i)}) - y_{obs}^{(i)})^2 + \lambda \sum_j \beta_j^2$$

larger weights penalize more due to squaring.

Lasso (L1):

$$\sum_i (y_i - \beta_0 - \sum_j \beta_j x_{ij})^2 + \lambda \sum_j |\beta_j| = \text{RSS} + \lambda \sum_j |\beta_j|$$

- λ proportional to β_j
- increasing λ increases bias, lowers variance
- perform feature selection as coefficients get to zero

comple

$$J(\beta_0, \beta_1) = \frac{1}{2m} \sum_i ((\beta_0 + \beta_1 x_{obs}^{(i)}) - y_{obs}^{(i)})^2 + \lambda \sum_j |\beta_j|$$

slower to converge than Ridge

Elastic Net:

$$J(\beta_0, \beta_1) = \frac{1}{2m} \sum_i ((\beta_0 + \beta_1 x_{obs}^{(i)}) - y_{obs}^{(i)})^2 + \lambda \sum_j (\alpha \beta_j^2 + (1-\alpha) |\beta_j|)$$

$\underbrace{\hspace{1cm}}$ $\underbrace{\hspace{1cm}}_{L1}$
~~Ridge~~ Lasso

RFE:

repeatedly measures feature importance, removes less important features.

$$\text{alphas} = []$$

$$\text{ridgeCV} = \text{RidgeCV}(\text{alphas} = \text{alphas}, \alpha=4). \text{fit}(\text{x_train}, \text{y_train})$$

$$\text{ridgeCV.alphas}, \text{ridgeCV.rmse} = \text{rmse}(\text{y_test}, \text{ridgeCV.predict}(\text{x_test}))$$

$$\text{LassoCV} = \text{LassoCV}(\text{alphas}, \text{max_iter}, \text{cv}). \text{fit}()$$

~~LassoCV.alphas~~
 LassoCV.rmse

$$\text{elasticNetCV} = \text{ElasticNetCV}(\text{alphas}, \text{l1_ratio} = [], \text{max_iter}). \text{fit}(,)$$

$$\text{elasticNetCV.alphas}, \text{.l1_ratio},$$

Classification

Supervised Learning

- Logistic Regression:

$$P(x) \frac{1}{1+e^{-(\beta_0 + \beta_1 x)}} = \frac{e^{(\beta_0 + \beta_1 x)}}{1+e^{(\beta_0 + \beta_1 x)}}$$

odds ratio, $\log\left(\frac{P(x)}{1-P(x)}\right) = \beta_0 + \beta_1 x$

(Ridge fit)

(R. fit(), penalty = 'l2', c = 10.0)

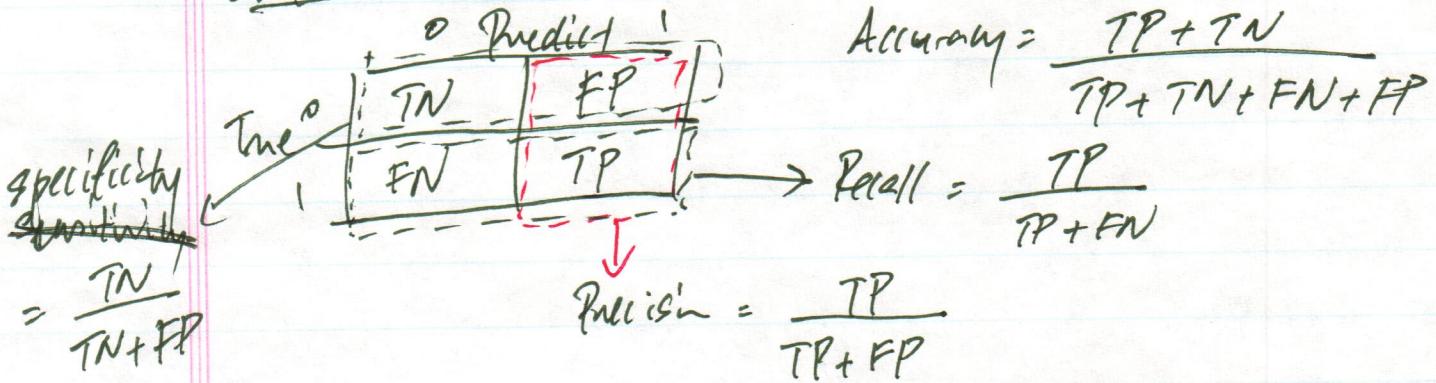
(R. fit(),)

(R. predict(),)

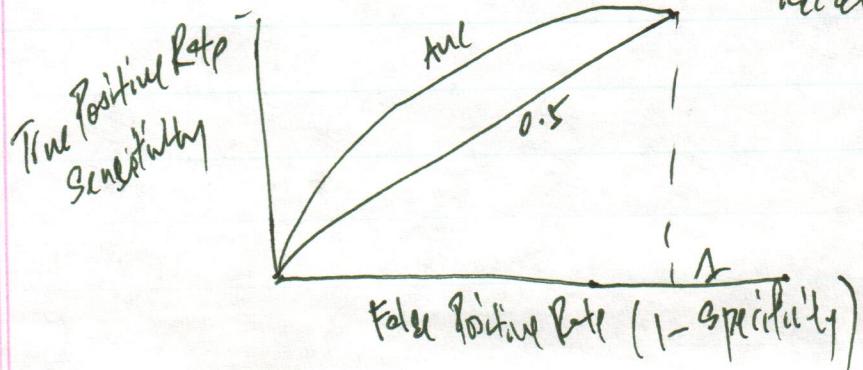
(R. coef(),)

Logistic Regression CV (penalty = 'l2', C = [1, 10, 100], cv = 4). fit(),)

Confusion Matrix:

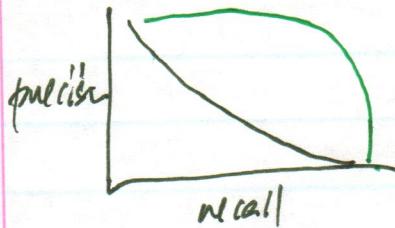


$$F1 = 2 \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$



Receiver operating characteristic (ROC)

Precision - Recall Curve - imbalanced dataset



accuracy_score(y-test, y-pred)

precision_score(,)

recall_score(,)

f1_score(,)

confusion_matrix

roc_curve, roc_auc_score

precision_recall_curve

lr_11 = LogisticRegressionCV(Cs=2] cv=4, penalty='l1',
solver='liblinear').fit(,)

lr_11.predict()

lr_11.predict_proba()

anc = roc_auc_score(label_binarize(y-test, classes=[0,1,2,3,4,5]),
y-proba, average='weighted')

cm = confusion_matrix(y-test, y-pred)

snf.heatmap(cm)

KNN:

manhattan distance $\sum_{x,y} |x-y|$ L1 euclidean distance $\sqrt{x^2 + y^2}$

KNN = kNeighborsClassifier(n_neighbours=3, similarity, weights)

KNN.fit(,)

KNN.predict()

high k, high bias, low k, high variance

categorical } LabelEncoder or pd.get_dummies(df, columns, drop-first=True)
 ordered } OrdinalEncoder
 binary } LabelBinarizer

StandScaler or MinMaxScaler

for scalar-fit-transform [df [column]]
 k in range(1, max_k):

KNN.fit(,)

y-pred = KNN.predict()

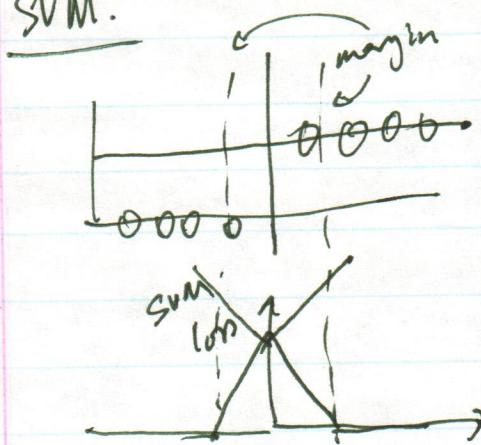
classification_report(y-test, y-pred)

accuracy_score(" ", " ")

f1-score(" ", " ")

cm(" ", " ")

SVM:



$$J(\beta_i) = \text{SVMCost}(\beta_i) + \frac{1}{C} \sum_i \beta_i$$

smoother C , more regularization,
more simple model

linearSVC < LinearSVC (penalty='l2', C=10)

LinearSVC.fit(,)

LinearSVC.predict()

lower β , C , more regularization, simpler model

rbfSVC = SVC(kernel='rbf', gamma=1.0, C=10)

rbfSVC.fit(,), rbfSVC.predict()

<u>Features</u>	<u>Data</u>	<u>Model choice</u>
many ($n > 10k$)	small (1k rows)	simple, logistic, linear SVC
few (< 100)	medium (10k rows)	SVC or RBF
few (< 100)	large ($> 100k$ rows)	odd features, logistic, linear SVC, kernel approx.

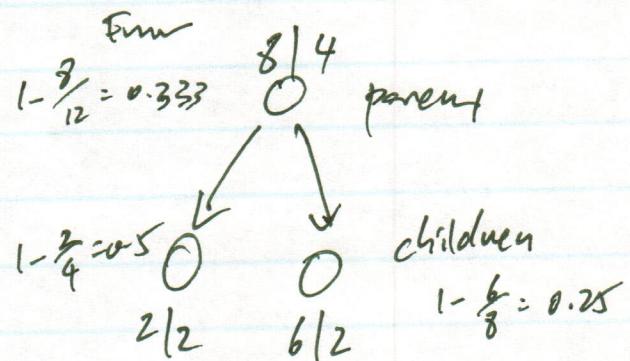
$\text{nystromSVC} = \text{Nyström SVC}(\text{kernel} = 'rbf', \text{gamma}, n\text{-components})$
 $X_{\text{train}} = \text{Nyström SVC}.fit_transform(X_{\text{train}})$
 $X_{\text{test}} = \text{Nyström SVC}.transform(X_{\text{test}})$
similar

$\text{rbfSampler} = \text{RBF Sampler}(\text{gamma}, n\text{-components})$

`sns.pairplot()`

Decision Tree:

- keep splitting until
- 1. only pure class remains
- 2. max. depth reached
- 3. performance metric reached



Classification error during split
 $E(t) = 1 - \max[p_{\text{left}}(t)]$

Classification error change

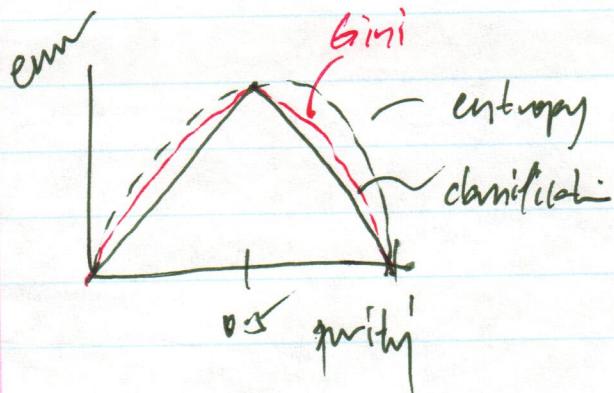
$$\begin{aligned}
 &= 0.333 - \frac{4}{12} \times 0.5 - \frac{8}{12} \times 0.25 \\
 &= p \\
 &\text{no further split}
 \end{aligned}$$

Split based on Entropy:

$$H(t) = - \sum_i p(i|t) \log_2 [p(i|t)]$$

$$\begin{array}{c} 8/4 \\ \swarrow \downarrow \\ 0 \quad b/2 \\ 2/2 \end{array} \quad -\frac{8}{12} \log_2 \left(\frac{8}{12}\right) - \frac{4}{12} \log_2 \left(\frac{4}{12}\right) = 0.9183$$
$$-\frac{6}{8} \log_2 \left(\frac{6}{8}\right) - \frac{2}{8} \log_2 \left(\frac{2}{8}\right) = 0.8113$$

$$\begin{aligned} \text{Classification error} &= 1 - \max [p(i|t)] \\ &= 0.9183 - 1 \times \frac{4}{12} = 0.8113 \times \frac{8}{12} \\ &= 0.0441 > 0, \text{ further split allowed} \end{aligned}$$



$$\frac{\text{Gini}}{h(t)} = 1 - \sum p(i|t)^2$$

Decision trees - high variance - overfit

- handle any data - binary, ordinal, categorical, numerical
- no scaling, no encoding

Regressor

dtc = DecisionTreeClassifier (criterion='gini', max_features=5,
max_depth=5)

dtc.fit(),
dtc.predict(),

dtc.tree_.nodecount, dtc.tree_.max_depth

```

param_grid = {
    'max_depth': range(1, dtc.tree_.max_depth + 1, 2),
    'max_features': range(1, len(dtc.feature_importances_) + 1, 2)
}

```

```

GR = GridSearchCV (DecisionTreeClassifier(random_state=42),
                    param_grid=param_grid,
                    scoring='accuracy',
                    n_jobs=-1)

```

GR.fit(X_train, Y_train)

GR.best_estimator_, tree_.node_count
 " " " " . max_depth

'neg_mean_squared_err'

```

export_graphviz(dtc, out_file='dot-data', filled=True)

```

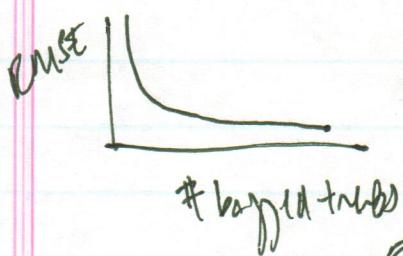
```

graph = pydotplus.graph_Len_dot_data(dot_data.getvalue())
graph.write_png('x.png')

```

Ensemble:

Bagging = bootstrap aggregating



more trees, less overfit

bagging trees compute in parallel.

bagged variance = $\frac{\sigma^2}{n}$ if trees independent

BC = Bagging Classifier(^{Regression}
 n_estimators = n)

BC.fit(,)

BC.predict()

Random Forest:

Random subset of features ← mechanism: $m/3$ + bootstrap

classification: J^m

Regression

`rf = RandomForestClassifier(n_estimators=4)`

`rf.fit(,)`

`rf.predict()`

Extra Random Trees:

Features + data randomly split - not greedy split.

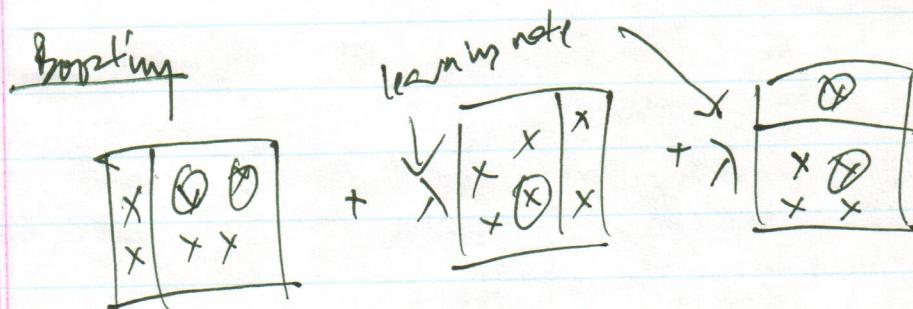
Regression

`etc = ExtraTreesClassifier(n_estimators=4)`

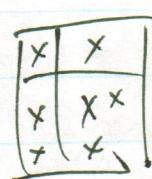
`etc.fit(,)`

`etc.predict()`

`rf.predict_proba()` → ^{rot-} acc-score
`rf.feature_importances_`



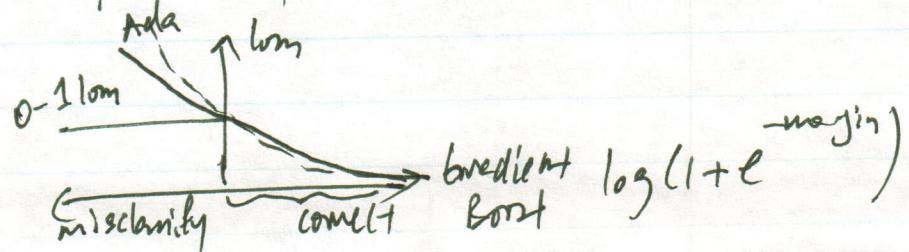
- initial weights
↓ wrong
- decrease weights
↓ right



more complex

more trees with lower learning rate, more regularization
 Higher learning rate leads to overfitting, less regularization

AdaBoost:

$$lm = e^{-margin}$$


Bagging

- bootstrap samples
- created independently / parallel
- only data points considered
- no weights
- will not overfit

Bagging

- entire dataset
- created successively
- residuals from previous models
- up-weight misclassified points
- reduce overfit

Tune GB model

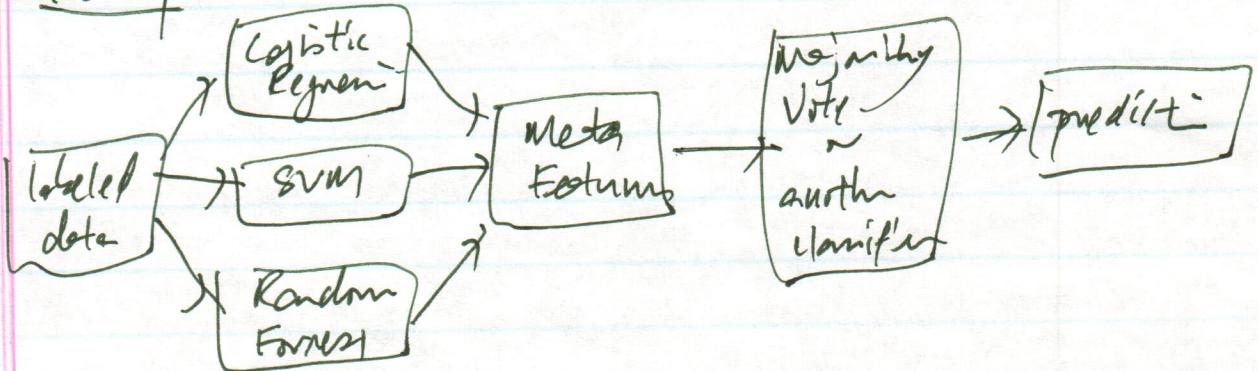
- $\lambda < 1$ for regularization
- subsample < 1.0 for fraction of data - SGD
- max_features - limit feature

`gbc = GradientBoostingClassifier(lr = 0.1, max_features = 1,
 subsample = 0.5, n_estimators = 200)`

`gbc.fit(,)`
`gbc.predict()`

`abc = AdaBoostClassifier(base_estimator = DecisionTreeClassifier(),
 n_estimators = n)`

Stacking:



class based on
prob. instead
of class
majority
vote

List of
fitted models

Regression

`vc = VotingClassifier(estimator_list, voting='soft')`

`vc.fit(,)`

`vc.predict()`

Stacking classifier / Regressor (`estimator_list`, `final_estimator = LogisticRegression()`)

Imbalanced Dataset:

Downsample - add more importance of minority

- \uparrow recall, \downarrow precision
- drop majority data.

Upsample - \uparrow recall, \downarrow precision

- fit duplicate minority class, - **Random Oversample** good for categorical

Resample - stratified split \rightarrow train test split

Shuffled split

kfold \rightarrow stratified fold

Synthetic Overampling

- choose 'K nearest neighbor' in minority class

- SMOTE - regular

- borderline - outlier, safe, in-danger

1 - connect in-danger to minority

2 - connect in-danger to whatever's nearby

- SVM - create new points from borders.

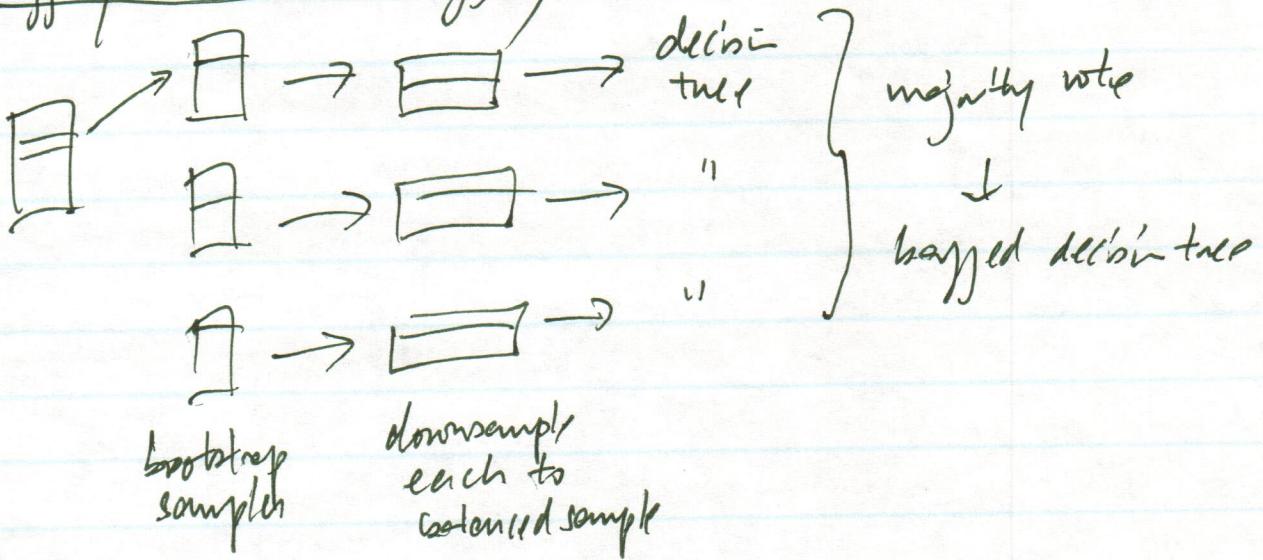
- ADASYN - move weight on values that might be misclassified.

by add more points away from minority area.

Undersampling -

- NearMin 1 - drop majority points closest to minority
 - throw off by outliers
- NearMin 2 - drop closest distance points in minority
- 3 -
- Tomek links
- Edited Nearest neighbor - run kNN, remove points don't \approx neig.

Bbagging - Balanced Bagging:



Always do train-test split first before under/oversample
Use sensible metric - F1, AUC, Cohen's Kappa

Unsupervised

clustering

- K-Means

- Hierarchical Agglomerative clustering

- DBSCAN

- mean shift

Dimensionality Reduce

- PCA

- Non-negative matrix factorization (NMF)

ANN	Multi-Layer Perception	Tutorial
RNN	RNN, LSTM	Predictive problems model sequences, forecasting
CNN	CNN	object recognit., segmentation in video/images
Unsupervised	Autoencoders, GANs	generate images, labeling dimension reduce.

Gradient Descent:

$$-\nabla J(\beta_0, \beta_1) \rightarrow \text{gradient of greatest decrease}$$

$$= \frac{\partial J}{\partial \beta_0}, \frac{\partial J}{\partial \beta_1}$$

$$w_i = w_o - \alpha \nabla \underbrace{\frac{1}{2N} \sum_i^N ((\beta_0 + \beta_1 x_{\text{obs}}^{(i)}) - y_{\text{obs}}^{(i)})^2}_{\text{learning rate cost funct.} = J}$$

next current learning rate cost funct. = J

until convergence / global minimum.

$$\text{gradient } \nabla J = \frac{\partial J}{\partial \beta}$$

$$= \frac{1}{N} \sum_i^N ((\beta_0 + \beta_1 x_{\text{obs}}^{(i)}) - y_{\text{obs}}^{(i)}) x_{\text{obs}}^{(i)}$$

Stochastic Gradient Descent

use only a single point

$$w_i = w_o - \alpha \nabla \frac{1}{2} ((\beta_0 + \beta_1 x_{\text{obs}}^{(i)}) - y_{\text{obs}}^{(i)})^2$$

more

~~noise~~ noise on only 1 point.

mini-batch:

$$n = 1 \text{ to } N$$

$$\omega_i = \omega_0 - \alpha \nabla \frac{1}{2n} \sum_i^n ((\beta_0 + \beta_i) \hat{y}_{obs}^{(i)}) - y_{obs}^{(i)})^2$$

- reduce memory
- less noisy than SGD

GD loop

for i in range (iter):

$$\hat{y} = \text{np.dot}(\theta^T \cdot X^T)$$

$$\text{loss_vec} = \text{np.sum}((y - \hat{y})^2)$$

$$\text{grad_vec} = (y - \hat{y}) \cdot \text{dot}(X) / n$$

$$\theta = \theta + lr * \text{grad_vec}$$

$$h(\theta) = \theta_0 + \theta_1 x$$

$$\text{cost function: } J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h(\theta)^i - y^i)^2$$

$$\text{gradient: } \nabla J = \frac{\partial J}{\partial \theta} = \frac{1}{m} \sum_{i=1}^m (h(\theta)^i - y^i) \cdot X_i^i$$

$$\text{update: } \theta_j = \theta_{j-1} - \alpha \frac{1}{m} \sum_{i=1}^m (h(\theta)^i - y^i) \cdot X_i^i$$

def cost(θ, X, y)

$$c = \frac{1}{len(y)} * \text{np.sum}(\text{np.square}((X \cdot \text{dot}(\theta)) - y))$$

return c

```
def gradient_descent(X, y, theta, alpha, iteration)
    initize : theta = np.zeros((iteration, 2))
    costs = np.zeros(iteration)
    m = len(y)
```

for i in range(iterations):

$$\theta = \theta - \frac{1}{m} * \alpha * (X^T \cdot \text{dot}(X \cdot \text{dot}(\theta)) - y))$$

$$\theta[i, :] = \theta^T$$

$$wts[i] = \text{cost}(\theta, X, y)$$

return ~~theta~~, theta, wts

Backpropagation:

imagine neural net as a function. $F: X \rightarrow Y$

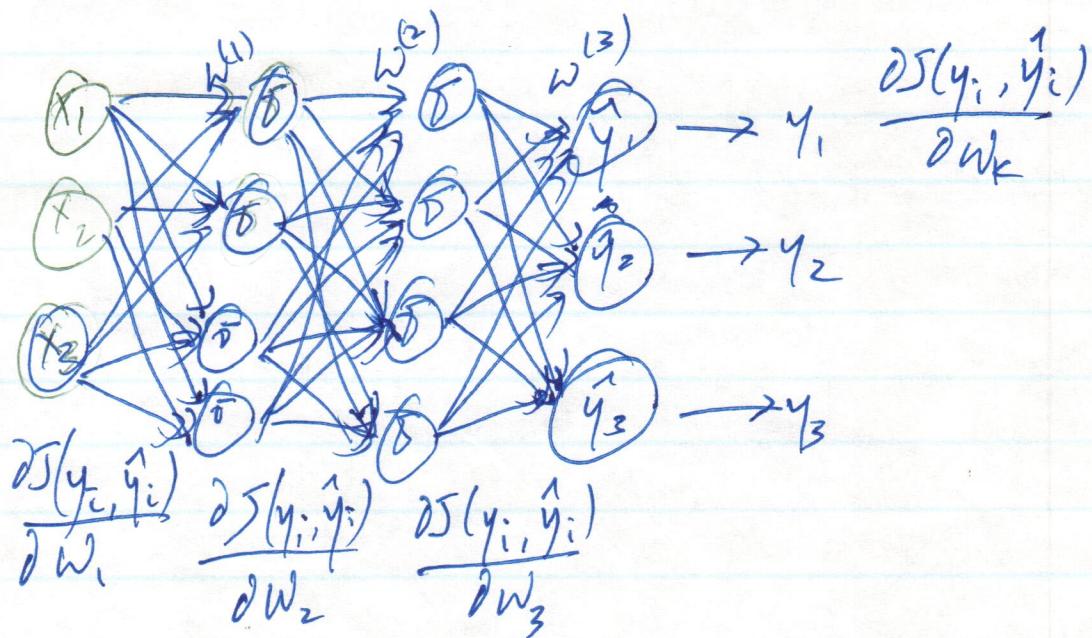
F computes many weights

(in function: $J(y, F(x))$)

$$F = \{ z_n = x_1 w_1 + x_2 w_2 + \dots + x_n w_n + b \}$$

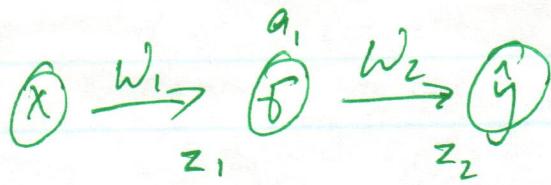
calculate $\frac{\partial J}{\partial w_k}$, update w_k

$$\frac{\partial J}{\partial w^{(3)}} = (\hat{y} - y) \cdot \cancel{a^{(3)}} \quad \frac{\partial J}{\partial w^{(2)}} = \frac{\partial J}{\partial w^{(3)}} \cdot a^{(2)} \quad \frac{\partial a^2 z^2}{\partial w} = \delta(z)(1 - \delta(z^2)) \\ \frac{\partial J}{\partial w^{(2)}} = (\hat{y} - y) \cdot \cancel{a^{(3)}} \cdot \cancel{a^{(1)}} \quad \frac{\partial J}{\partial w^{(1)}} = (\hat{y} - y) \cdot w^{(2)} \cdot \delta'(z^{(2)}) \cdot \cancel{a^{(1)}} \cdot \cancel{a^{(0)}} \cdot \cancel{x}$$



$$\delta'(z) = \delta(z)(1 - \delta(z)) \leq 0.25 \quad \text{for sigmoid, max is 0 to 1, so max } 0.5(1 - 0.5) = 0.25$$

more layers, gradient gets smaller, vanishing gradient.



def forward_pass(w_1, w_2):

$$z_1 = \text{np.dot}(X, w_1)$$

$$a_1 = \sigma(z_1)$$

$$z_2 = \text{np.dot}(\cancel{a_1}, w_2)$$

$$\hat{y} = \sigma(z_2)$$

$$\frac{\partial J}{\partial z_2} = \hat{y} - y$$

$$\frac{\partial J}{\partial w_2} = \text{np.dot}\left(\frac{\partial J}{\partial z_2}, a_1\right) = \boxed{(\hat{y} - y)a_1}$$

$$\frac{\partial a_1}{\partial z_1} = \sigma'(z_1) = \sigma(z_1)(1 - \sigma(z_1))$$

$$\frac{\partial J}{\partial w_1} = \text{np.dot}\left(\frac{\partial J}{\partial z_2}, w_1\right) \sigma'(z_1) \times = \boxed{(\hat{y} - y)a_1} w_1, \sigma'(z_1) \times$$

$$\text{return } \left(\frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial w_2} \right), \hat{y}$$

def loss_fn(\hat{y}, y_{true} , eps=1e-6):

$$\text{loss} = -\frac{1}{N} \left(\hat{y}_{\text{true}} \log(\hat{y}) + (1 - \hat{y}_{\text{true}}) \log(1 - \hat{y}) \right)$$

initial w_1, w_2

for i in range(num_iter):

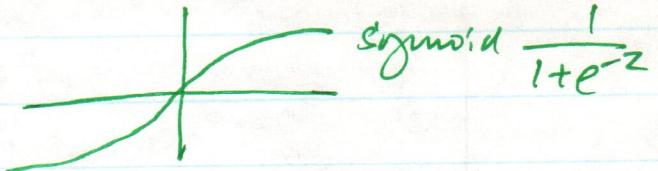
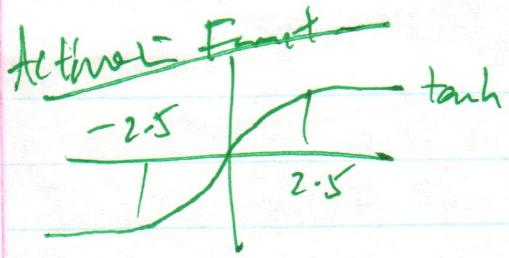
$$\hat{y}, \left(\frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial w_2} \right) = \text{forward_pass}(w_1, w_2)$$

$$w_1 = w_1 - lr \frac{\partial J}{\partial w_1}$$

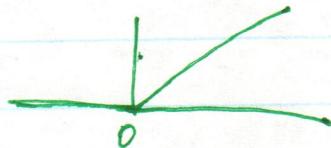
$$w_2 = w_2 - lr \frac{\partial J}{\partial w_2}$$

$$\text{loss} = \text{loss_fn}(\hat{y}, y)$$

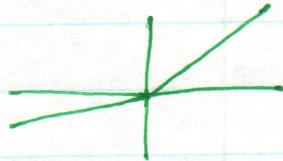
$$\text{acc} = \text{np.sum}(\hat{y} > 0.5 == y) / N$$



$$\text{ReLU}(z) = \begin{cases} 0 & z < 0 \\ z & z \geq 0 \end{cases}$$



$$\text{Leaky ReLU}(z) = \begin{cases} \alpha z & z < 0 \\ z & z \geq 0 \end{cases}$$



Regularization

- regularization penalty in loss function $J = \frac{1}{2n} \sum_i^n (y^i - \hat{y})^2 + \lambda \sum_i w_i^2$
- dropout - randomly remove subsets of neurons.
- early stopping - no change in validation after patience of n
- stochastic / mini-batch GD

ridge

nn

w

Optimizers

$$\omega := \omega - \alpha \nabla J$$

momentum - $v_t = \eta v_{t-1} - \alpha \cdot \nabla J$

(layer steps) $\omega := \omega - v_t$

if moving in

~~same~~ some direction.

Vestor momentum -

Adagrad - scale update for each weight

$$\omega := \omega - \frac{1}{\sqrt{\theta}} \nabla J$$

WT

RMSProp - decay older gradients more than currents.

Adam - combines momentum + RMSProp

Epoch - one pass through all training data.

Full batch GD - 1 step per epoch

SGD - n steps per epoch ($n = \text{training size}$)

minibatch - $\frac{n}{\text{batchsize}}$ steps per epoch

Keras:

```
model = Sequential()
```

```
model.add(Dense(units=4, input_shape=[3]))
```

```
model.add(Activation('sigmoid'))
```

```
model.add(Dense(units=4))
```

```
model.add(Activation('sigmoid'))
```

```
model.summary()
```

```
optimizer='adam', loss='binary_crossentropy'
```

```
model.compile(optimizer='sgd', loss='binary_crossentropy',
```

```
metrics=['accuracy'])
```

```
history = model.fit(X_train, y_train,  
                     validation_data=(X_test, y_test),  
                     epochs=n, shuffle=True)
```

```
y_pred_class = model.predict_classes(X_test)
```

```
y_pred_proba = model.predict(X_test)
```

~~```
y_pred_class = np.argmax(y_pred_class, axis=1)
```~~

Multiclass classification:

$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum_k e^{z_k}}$$

$$\text{categorical cross entropy} = -\sum_i y_i \log(\hat{y}_i)$$

$$\frac{\partial \text{CE}}{\partial \text{softmax}} \cdot \frac{\partial \text{softmax}}{\partial z_i} = \hat{y}_i - y_i$$

Scaling:

min max  $(0, 1)$ ,  $x_i' = \frac{x_i - x_{\min}}{x_{\max} - x_{\min}}$

linear scaling  $(-1, 1)$ ,  $x_i' = 2\left(\frac{x_i - x_{\min}}{x_{\max} - x_{\min}}\right) - 1$

Convolution:

Padding

Stride

convolution size

$$\text{output size} = \frac{(\text{input size} + 2 * \text{padding} - (\text{kernel size} - 1))}{\text{stride}} + 1$$

Pooling - reduce image size  
- maxpooling, avg pooling

Transfer Learning

Convolution  $\rightarrow$  Fully Connected layer  $\rightarrow$  Predict  
 $\rightarrow$  new layers  $\rightarrow$  new Predict

```

if K.image_data_format() == 'channel_last':
 input_shape = (1, img_rows, img_cols)
else:
 input_shape = (img_rows, img_cols, 1)

```

frozen layers:

```

for l in feature_layers:
 l.trainable = False

```

CNN:

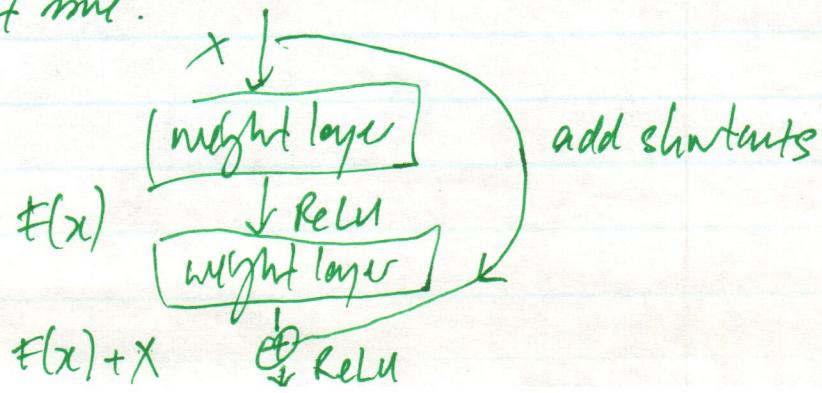
LeNet

AlexNet - data augmentation - crop, flipping, zoom  
 VGG - deep network w 3x3 convolutions.

Inception -



ResNet - early layers slow to adjust due to vanishing gradient issue.



RNN:

Forecasting, speech, genome,

LSTM:

Forecasting, speech, anomaly detect.

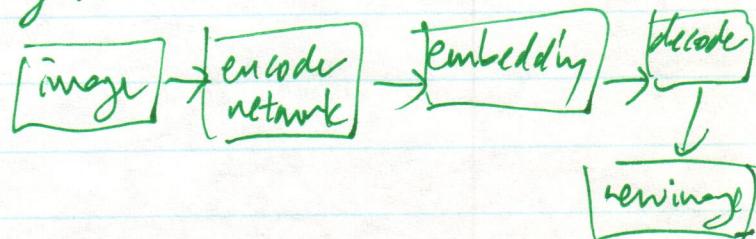
Autoencoders:

- learn similarity / differences in images.

- reduce dimensionality of images.

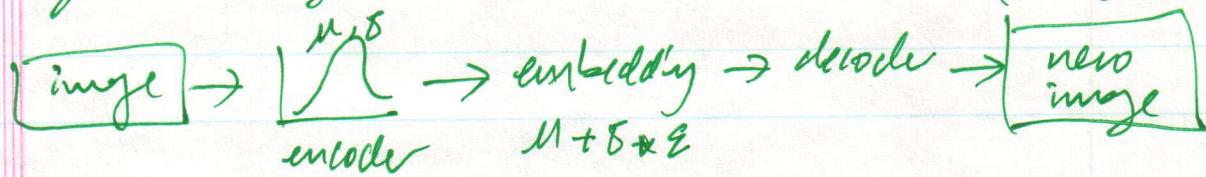
{ anomaly detect -

day discovery - ...



Variational AE:

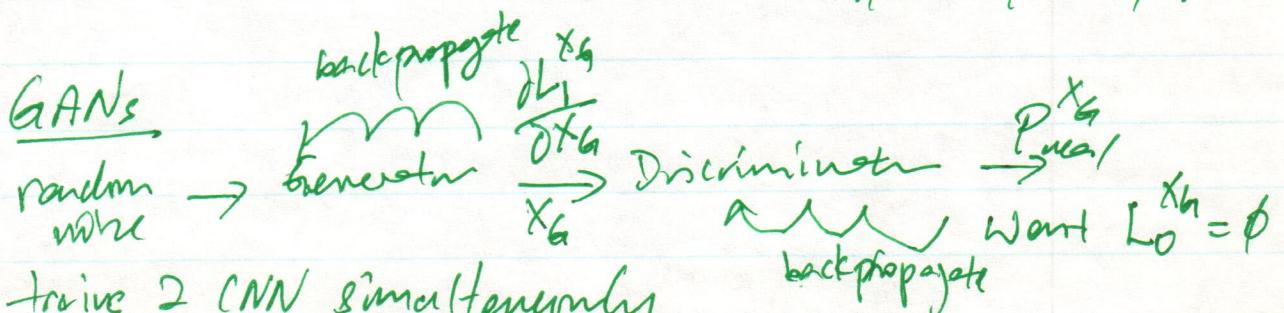
generate images base on normal distribution of images.



penalty for not constructing new images correctly

generating  $M, B$  equal to 0 and 1

$$\text{KL Divergence} = \frac{1}{2} \left( \underbrace{e^{\log(\theta)} (\log(\theta) + 1)}_{\text{penalizes for } \log(\theta) \neq M} + \underbrace{M^2}_{\text{different from } \theta} \right)$$



training 2 CNN simultaneously

$P_g$  - probability generated image is real

compute (on assuming  $P_g = \emptyset$  (false))

$$L_0^{X_k} = f(P_{\text{real}}^{X_k}, \phi)$$

$$L_1^{X_k} = f(P_{\text{real}}^{X_k}, 1)$$

More sensitive to hyperparameters than other DL models.

CIME (Closely-Interpretable Model Agnostic Explanations):

- sensitivity analysis of outputs to inputs -

### Reinforcement Learning

Agents interact in Environment

- choice of Actions in response to state limited by Policy
- impacts by Rewards, not known.