

Types of Machine Learning

- (1) Supervised: $f(x_i) = y_i$ (most success)
e.g. spam detect, medical diag, click pred
- (2) Unsupervised: x_i comes from $p(x)$ & we want to learn p (hard to evaluate)
e.g. clustering, dim red, topic modeling
- (3) Reinforcement Learning: neural nets, agents interact w/ envir & rewards
e.g. games, robotics, self-driving
- (4) Other: semi-supervised, active learning, forecasting

Data Visualization

- (1) Show data & maximize data-ink ratio
- (2) Tools matter
- (3) Labels - title, axes, color, legend, etc.

Channels for quantitative data

(1) Position	(5) Volume	* Think hard about aspect ratios, baselines, & 3D plotting
(2) Length	(6) Texture density	
(3) Angle	(7) Color saturation	
(4) Area	(8) Color hue	

fig, ax = plt.subplots(n, m)
 $ax[0, 0].plot(x, y)$
 $ax[0, 1].scatter(x, y, alpha=0.1)$
 $ax[1, 0].hist(x, bins=1000)$
 $ax[1, 1].bar(range(len(y)), y)$

Heatmap: color code to represent diff values

Hexgrid: density map (brighter colors = more pts)

Nearest Neighbors * $k \rightarrow \infty$, overfitting!
 Assign data label of avg of nearest neighbors

KNN = KNeighborsClassifier(n_neighbors=3)

knn.fit(x-train, y-train) * fast to train
 $y_pred = knn.predict(x_test)$ * large to store
 * slow to predict

Regression $\hat{y} = W^T x + b = \sum_i w_i x_i + b$

DLS: minimize $\sum_{i=1}^n \| \hat{y}_i - y_i \|^2$ * unique soln if X full rank

Ridge: $\min_{W, b} \sum_i (\hat{y}_i - y_i)^2 + \alpha \| W \|^2$
 Always has soln if $\alpha > 0$
 high α = underfit
 small α = overfit

Accur



Lasso: $\min_{W, b} \sum_i (\hat{y}_i - y_i)^2 + \alpha \| W \|^1$
 Sets some w_i to 0 = auto feature select

Elastic Net: $\min_{W, b} \sum_i (\hat{y}_i - y_i)^2 + \alpha \| W \|^1 + \beta \| b \|^2$
 2 hyperparams

Coeff of Determination: $R^2(y, \hat{y}) = 1 - \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \bar{y})^2}$
 * can be neg for biased data / test set

* For skewed data, log transform better results

Classification

Binary case: $\hat{y} = \text{sign}(W^T x + b)$
 * gives prob estimates

Logistic Regression: Loss $\rightarrow \min_{W, b} \sum_i \log(\hat{y}_i / y_i) + (1 - \hat{y}_i)^2$

$\log(\frac{p(y=1|x)}{p(y=0|x)}) = W^T x + b \rightarrow p(y=1|x) = \frac{1}{1 + e^{-W^T x - b}}$

Loss function: $\sum_i \log(\exp(-y_i(W^T x_i + b)) + 1)$

Add a penalty: $\min_{W, b} C \sum_i \log(\exp(-y_i(W^T x_i + b)) + 1) + \frac{1}{2} \| W \|^2$

$\| W \|^2$ \rightarrow big C less regular, small C more regular

SVM: draw boundary b/w classes w/ max margin (points kept outside margin, nearest pts = SV's)

Loss function: $\min_{W, b} C \sum_i \max(0, 1 - y_i(W^T x_i + b)) + \frac{1}{2} \| W \|^2$

W within margin $\rightarrow y_i(W^T x_i + b) < 1$ * big c smaller $W \rightarrow$ larger margin

Linear SVM: soft margin, $\| W \|^2$ or $\| W \|^2 + \text{margin loss}$

Kernel SVM: non-linear, nonlinear models

$\hat{y}(x_i) = \text{sign}(\sum_j \alpha_j y_j K(x_i, x_j)) \rightarrow \hat{y} = \text{sign}(\sum_j \alpha_j y_j W^T x_i)$

Multiclass Classification: * classify pt by all classifiers or return most common class

* softmax loss: $L(x, y) = -\sum_i y_i \log(\sum_j e^{W_j^T x_i})$

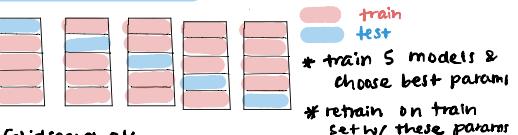
Softmax Loss: $L(x, y) = -\sum_i y_i \log(\sum_j e^{W_j^T x_i})$

$\min_{W, b} -\sum_i y_i \log(p(y_i | x_i, W, b))$

For high cardinality categories \rightarrow instead of OHE, one "response encoded" var (i.e. avg \$ for zip)

Cross-Validation

k-fold, k=5 or 10 usually



* train 5 models & choose best params

* retain on train set w/ these params

GridsearchCV param_grid = {n_neighbors: np.arange(1, 15, 2)}

grid = GridSearchCV(KNeighborsClassifier(), param_grid=param_grid, cv=10, return_train_score=True)

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

grid.best_score_

grid.best_params_

grid.score(X-test, y-test)

Different Strategies

Kfold: k models that split diff train & test * BAD for ordered data \rightarrow shuffle!

Stratified Kfold: Take some % of data for test & train across diff classes * good for imbalanced classes

Leave-one-out: k-fold ($n_folds = n_samples$) * high variance, takes too long

2 shuffle Split: repeat sample test set w/ replace

2 Repeated KFold: repeatedly shuffle & run kfold

Group KFold: each group (i.e. city) should be entirely in train or test set

TimeSeriesSplit: "rolling window" of test/train (train should always be past of test)

Preprocessing / Pipelines

Scaling (esp imp for dist-based algorithms)

* Standard Scaler: subtract mean to get var

MinMax Scaler: min val=0, max val=1

Robust Scaler: uses medians & quantiles

Normalizer: divide each row by Euclid norm

* For sparse data, scale but don't center

s = StandardScaler()

s.fit_transform(X-train, y-train)

s.transform(X-test, y-test) * don't fit on test!

* MUST include preprocessing inside CV so we don't "leak info" and fit scaling on test data

pipe = make_pipeline(StandardScaler(), Ridge())

cross_val_score(pipe, X, y)

param_grid = {'ridge_alpha': [0.01, 0.1, ..., 1000]}

grid = GridSearchCV(pipe, param_grid=param_grid, cv=10, return_train_score=True)

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

Missing Data (99%, ?, np.NaN, "N/A", None...)

* sometimes missingness informative

Imputation: (1) mean/median (2) KNN (3) regression (4) Matrix Factorization

imp = SimpleImputer(strategy="mean").fit(X-train)

X-trans = imp.transform(X-train)

knn-pipe = make_pipeline(KNNImputer(), StandardScaler(), LogisticRegression())

(3) Regression: train regression models to predict missing col

* Retrain after imputing missing data

* for catg data, add "cat_missing" category

Categorical Vars

Ordinal Encoding: Brooklyn \rightarrow 1, Manhattan \rightarrow 2, Queens \rightarrow 3

last col = 1 - others can create collinearity

OHE (drop 1 + which impt for d models)

pd.get_dummies(df, columns=['Boro'])

OneHotEncoder().fit(df) \rightarrow transforms non-cat cols

Column Transformer

cat = df.dtypes == object

preprocess = make_column_transformer(

(StandardScaler(), ~cat),

(OneHotEncoder(), cat))

Model = make_pipeline(preprocess, LogisticRegression())

Target / Impact Encoding

For high cardinality categories \rightarrow instead of OHE, one "response encoded" var (i.e. avg \$ for zip)

Trees, Forests, Ensembles

Decision Tree: split based on threshold of feature w_i most into gain overall classes

Classification criteria: Gini = $\sum_k P(k)(1-P(k))$ prop of class k can't extrapolate, fast though, no preprocess

Cross-Ent = $-\sum_k P(k) \log P(k)$ in node m

* Regression trees exist $\rightarrow \hat{y}_m = \frac{1}{n} \sum_i y_i$ (MSE, MAE)

To prevent overfit: (1) prepruning (limit size while building)
 (2) postpruning (cut branches after build)

Cost complexity: $R(T) + \alpha |T| \rightarrow |T| = \# \text{leaf nodes}$ * BEST = greedy build tree then postprune w/ cost complexity

Probability = fraction class in leaf (without reg, no imports)

Ensemble: Build many models & avg (more models better)

Bagging: build diff models by bootstrapping data each usually uses 2/3 data

Extremely Random Tree: randomly choose thresh for feature, no bootstrap, fast!

Gradient Descent & Boosting

learning rate To minimize F(w), keep updating: $w_i \leftarrow w_i - \eta_i \frac{\partial F(w)}{\partial w_i}$

Stochastic GD: pick one i , compute partial gradient, & update w_i * fast for large datasets

Boosting: "weak" learners ensemble to create "strong" one $y = f_1(x) + f_2(x) + f_3(x) + \dots$ * # in learning rate δ $y = y + \delta f_2(x) + \delta f_3(x) + \dots$

Classification min $\sum_i \log(\exp(-y_i \hat{y}_i) + 1)$ Multiclass $\text{ply}(c|x) = \frac{e^{\hat{y}_c}}{\sum_j e^{\hat{y}_j}}$ Log Loss

* too many trees = overfitting (stop adding trees when val acc stagnates)

(1) pick # trees & learning rate (max depth good)
 (2) pick learning rate & use early stopping

XGBoost speeds up tree building by binning & columns subampling

Model Evaluation

Binary classification

Actual	Pred -	Pred +	TP	FP	FN	TN	Acc = TP + TN
-	True -	False +					TP + TN = FP + FN
+	False -	True +					# bad for class imbalance

Precision: $\frac{TP}{TP + FP}$ (+ pred value) Recall: $\frac{TP}{TP + FN}$ (sensitivity, true + rate)

F2: $\frac{2 \cdot \text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$ * generally care about TPR, FPR, precision

Averaging: Macro: $\frac{1}{k} \sum_i R(y_i, \hat{y}_i)$ treats all classes =
 Weighted: $\frac{1}{\sum_i n_i} \sum_i n_i R(y_i, \hat{y}_i)$ weights larger classes more

Balanced: $\frac{1}{k} \left(\frac{TP}{TP + FP} + \frac{TN}{TN + FN} \right)$

$y_pred = rf.predict_proba(X_test)[:, 1] > 0.50$

Precision-Recall curve: gives prec/recall over diff thresh

ROC curve: TPR vs. FPR * AUC = area under ROC (0.5 = random)

Avg precision: $\frac{1}{k} \sum_i P(c) \cdot \text{Recall}$ sum over data pts k ranked by decision fn

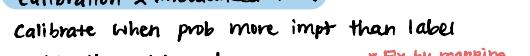
Regression: R², MSE, MAE (sensitive to scaling & outliers)

MAPE: $\frac{1000}{n} \sum_i \frac{|y_i - \hat{y}_i|}{y_i}$

Calibration & Imbalanced Data

Calibrate when prob more impt than label

calibration plot + (reliability diag)



* Fix by mapping classifier probs to better probs (f.calib)(x) = p(y)

* use w/ val set, not train set, not SVMs

Brier Score (bin class) = $\frac{1}{n} \sum_i (p(y_i) - y_i)^2$

Platt scaling: $f_{\text{platt}} = \frac{1}{1 + \exp(-w^T x - b)}$ good for SVMs

Isotonic Reg: flexible func (learns increasing step fn)

* optimum monotonic func on training data

Imbalanced Data (asymmetric cost or data)

Random undersampling, random oversampling

* RandomUnderSampler() RandomOverSampler()

Instead of resampling, reweight loss func

min $-C \sum_i y_i \log(\exp(-y_i(W^T x_i + b)) + 1) + \|W\|^2$

* SMOTE good sometimes Hinge = $\sum_i C_k P(k)(1 - P(k))$

Ensemble resampling: random resampling for each tree

Linear model coeffs useless w/o scaling, w/ regularization, diff feature correlations

Prop Feature import: $I_i^{prop} = \text{Acc}(x_i, \text{y}) - \text{Acc}(x_i, \text{x}_{-i}, \text{y})$ same model

Permutation import: $I_i^{perm} = \text{Acc}(x_i, \text{y}) - \text{Acc}(\text{x}_{-i}[\text{acc}(x_i, \text{x}_{-i}, \text{y})], \text{y})$

LIME: sparse lin local model around each datapoint

SHAP: dropout import for every subset of features

Partial Dep Plots: $f_i^{PDP}(x_i) = E_{\text{X}_{-i}}[f(x_i, \text{x}_{-i})]$ marginal dep of pred