

REGRESSION AND PREDICTION

- Lucy Xu



QUIZ

- ✗ What is Simple Linear Regression(MLR) assumptions?
- ✗ What is Cross Validation and why/how to do it?
- ✗ What's the logistic regression loss function?
- ✗ What's gradient descent?
- ✗ What's regularization?



TYPES OF MACHINE LEARNING

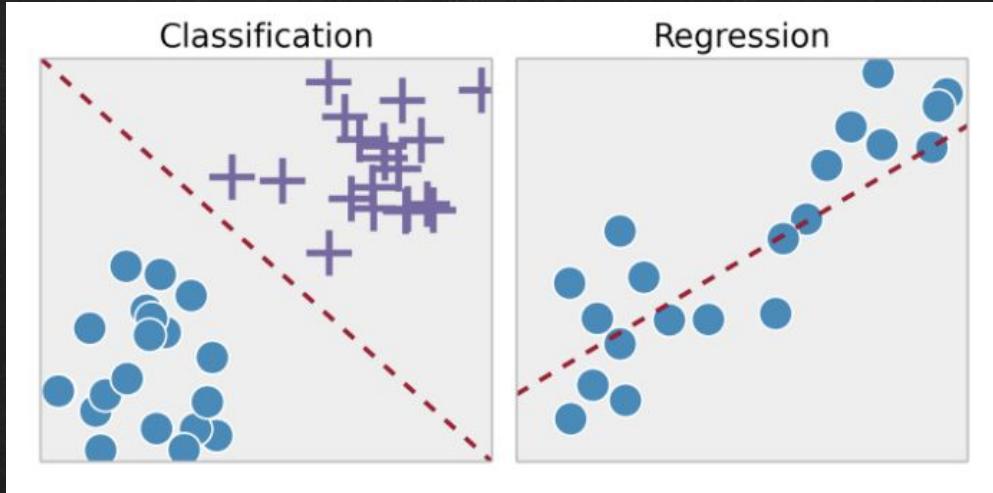
✗ Supervised Learning

- Target type: Continuous and Classification
- Algorithms: Linear regression, logistic regression, Decision Tree, Random Forest, XGB, Naive Bayesian, SVM, Neural Network, KNN

✗ Unsupervised Learning

- Algorithms: K-means, PCA

TYPICAL PROBLEMS OF SUPERVISED LEARNING



Decision: Approve/Decline,
Solicitation/not, market/not,
promotion/not;
Attrition: Leave/Stay,
happy/unhappy
Risk: Yes/No

House Value,
Product price,
Customer income

LINEAR MODELS

A model which is linear in the parameter

SLR model: $y = \beta_0 + \beta_1 x + \varepsilon$

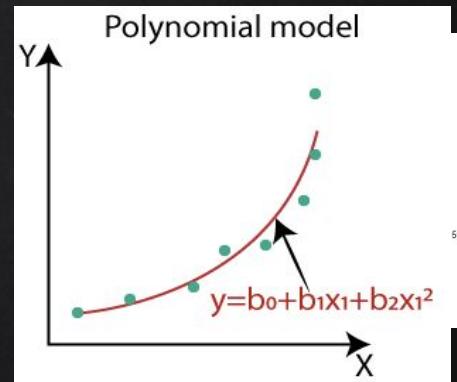
MLR model: $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$

Polynomial regression model:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_{1,\text{square}} + \varepsilon$$

Nonlinear model example:

Logistic growth model: $y = \alpha/\beta e(-kx) + \varepsilon$





TOPICS TO COVERED IN REGRESSION

- ✗ Simple linear regression
- ✗ Linear algebra 101
- ✗ Multiple-regression
- ✗ Logistic regression
- ✗ Gradient descent
- ✗ Regularization



1.

SIMPLE LINEAR REGRESSION

"THIS ALGORITHM ASSUMES THAT THERE IS A LINEAR RELATIONSHIP BETWEEN THE 2 VARIABLES, INPUT (X) AND OUTPUT (Y), OF THE DATA IT HAS LEARNT FROM. THE INPUT VARIABLE IS CALLED THE *INDEPENDENT VARIABLE* AND THE OUTPUT VARIABLE IS CALLED THE *DEPENDENT VARIABLE*. WHEN UNSEEN DATA IS PASSED TO THE ALGORITHM, IT USES THE FUNCTION, CALCULATES AND MAPS THE INPUT TO A CONTINUOUS VALUE FOR THE OUTPUT."

FOR EXAMPLE: IF X INCREASE Y WILL INCREASE TO A CERTAIN MAGNITUDE OR DECREASE TO A CERTAIN MAGNITUDE.



LINEAR REGRESSION KEY TERMS

Dependent variable – The variable is to predict. Also called: Response, Y variable, target, outcome

Independent variable – The variable used to predict the response. Also called: X variable, feature, attribute, predictor

Observation – The vector of predictor and outcome values for a specific individual or case. Also called, row, case, instance, example, record, $(x_1, y_1), (x_2, y_2)$

Intercept – The intercept of a regression line is the predicted value when $X=0$. Also called b_0 , β_0 ,

Regression coefficient – The slope of a regression line. Also called, slope, b_1 , β_1 , w_1

Fitted values – The predicts obtained from the regression line. Also called, estimates, predicted values, \hat{Y}_i

Residuals – The difference between the observed values and the fitted values. Also called, errors

Least squares – The method of fitting a regression by minimizing the sum of squared residuals. Also called, ordinary least squares(OLS, MSE)

Maximum Likelihood Estimation – “a method of estimating the parameters of a probability distribution by maximizing a likelihood function”, also called, MLE



SIMPLE LINEAR REGRESSION EQUATION

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

$$\hat{Y} = b_0 + b_1 X$$

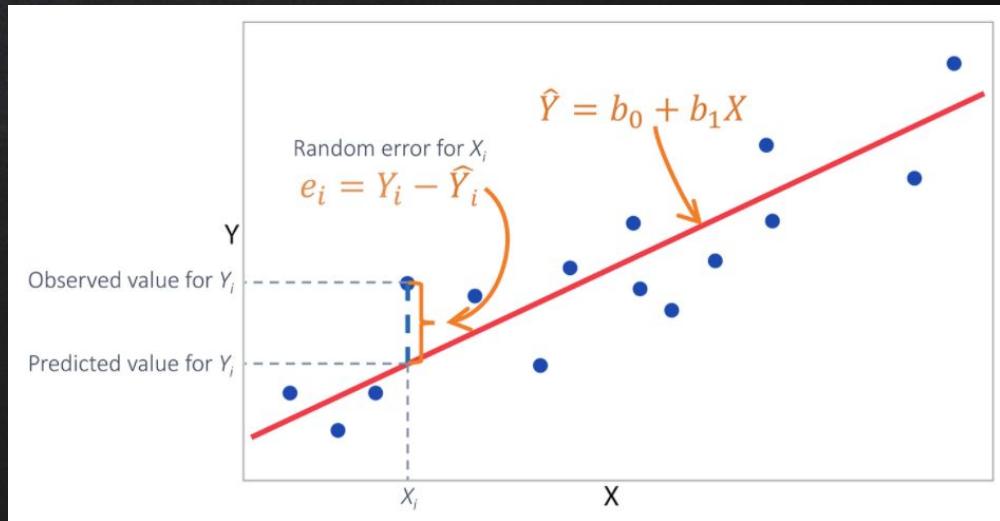
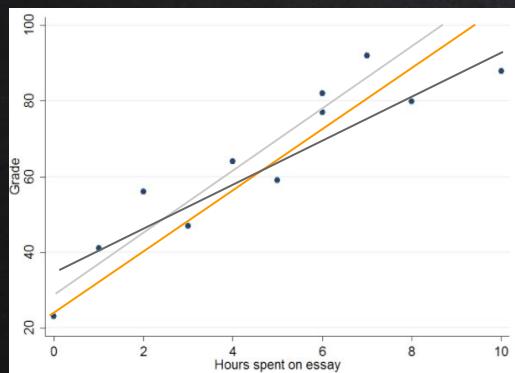
Simple Linear regression assumption:

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, i=1,2,3,4,\dots$$

1. Error is zero on average
2. Error has constant variance
3. Each observation is independent to other
4. ε_i are normal distributed

LEAST SQUARES ESTIMATION

With data $(x_1, y_1), (x_2, y_2), (x_3, y_3), \dots$, there may be several ways to fit a linear regression model





LEAST SQUARES ESTIMATION

With data $(x_1, y_1), (x_2, y_2), (x_3, y_3), \dots$, find the candidate model ($\hat{Y} = b_0 + b_1 X$) which minimize:

$$S = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - (b_0 + b_1 x_i))^2.$$

By mathematical deviation to find the value of b_0 and b_1 to make S the smallest of all candidates:

Intercept: $b_0 = \bar{y} - b_1 x_mean$

$$\text{Slope: } b_1 = r \frac{s_y}{s_x}$$

$$\text{Correlation Coefficient: } r = \frac{1}{n-1} \sum \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right)$$

OLS DERIVATION DETAIL



$$S(\beta_0, \beta_1) = \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

β_0 and β_1 are random variables; b_0, b_1 are the final values to minimize the square of error

The partial derivatives of $S(\beta_0, \beta_1)$ with respect to β_0

$$\frac{\partial S(\beta_0, \beta_1)}{\partial \beta_0} = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)$$

The partial derivatives of $S(\beta_0, \beta_1)$ with respect to β_1

$$\frac{\partial S(\beta_0, \beta_1)}{\partial \beta_1} = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) x_i$$

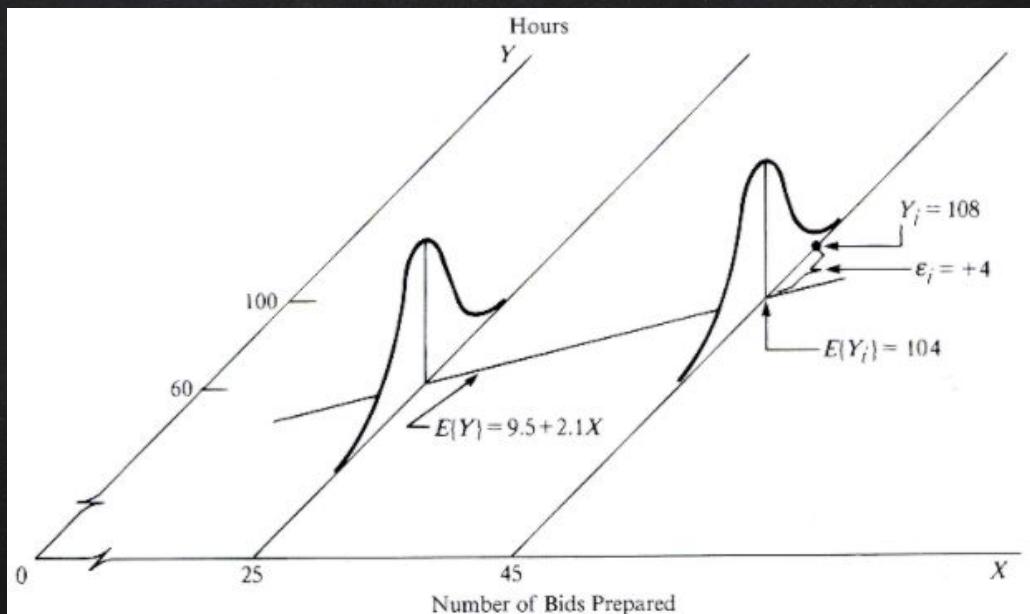
The solution (b_0, b_1) of β_0, β_1 when S is smallest are obtaining by setting,

$$\frac{\partial S(\beta_0, \beta_1)}{\partial \beta_0} = 0$$

$$\frac{\partial S(\beta_0, \beta_1)}{\partial \beta_1} = 0.$$

Or through Gradient Descent

REVISIT THE NORMAL ASSUMPTION

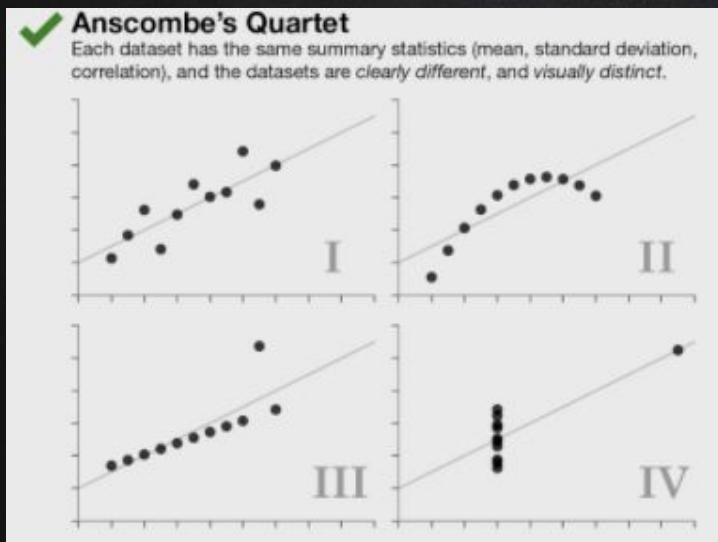


T test in LR: The coefficient b_1 for a predictor, divided by the standard error of the coefficient b_1 , giving a metric to compare the importance of variables in the model.



RESIDUAL DIAGNOSTIC

Regression use for prediction and explanation, your model should be intuitively reasonable beside with the mean error square (MSE)

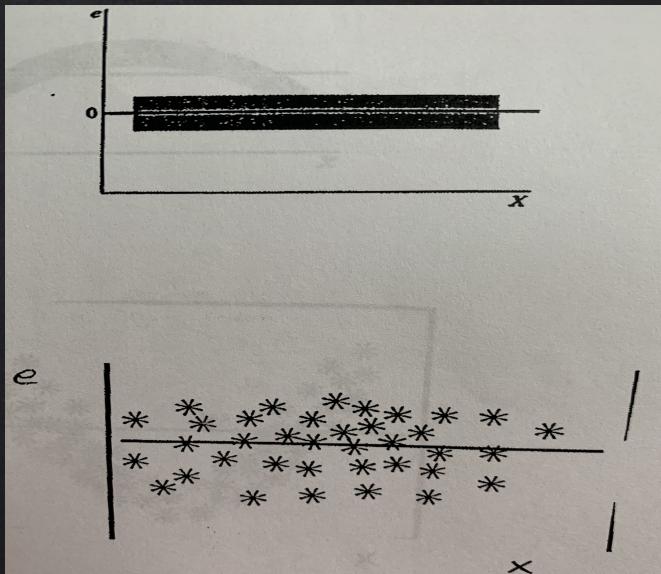


Always conduct a scatter plot for regression

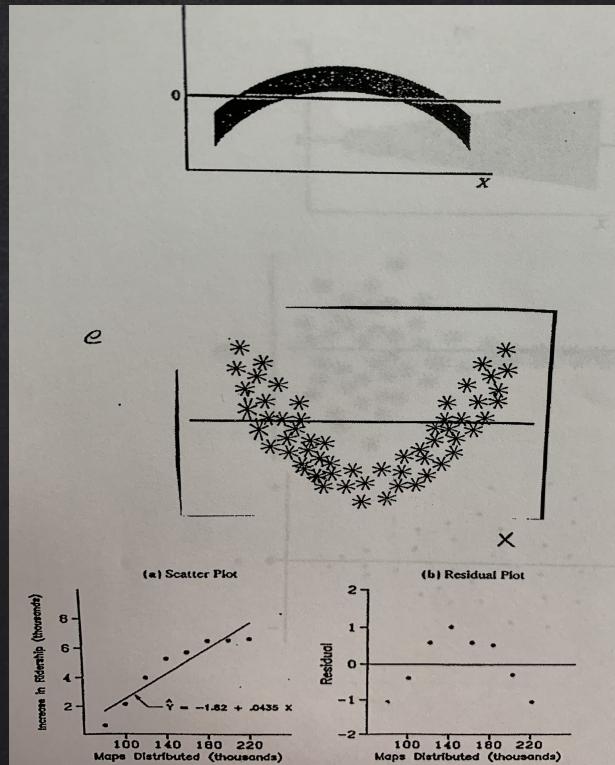
Residual diagnostic help disclose wrong inception and adjust business understand

RESIDUAL DIAGNOSTIC - IDEA RESIDUALE

$E_i = Y_i - \hat{Y}_i$, idea residual plot is like this:

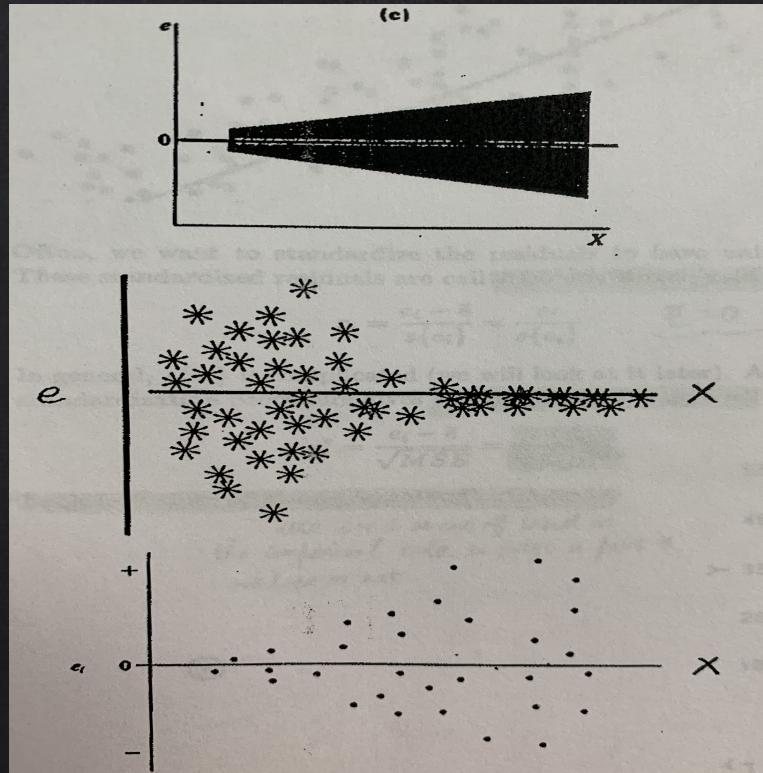


RESIDUAL DIAGNOSTIC - TRENDED RESIDUALS



~~D
S
L
X~~

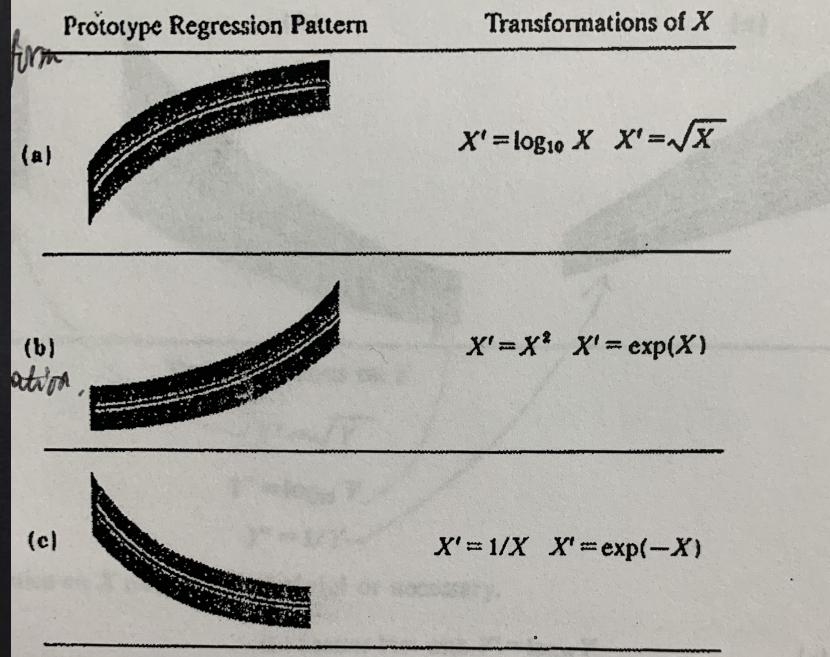
RESIDUAL DIAGNOSTIC - INCONSISTENT VARIANCE



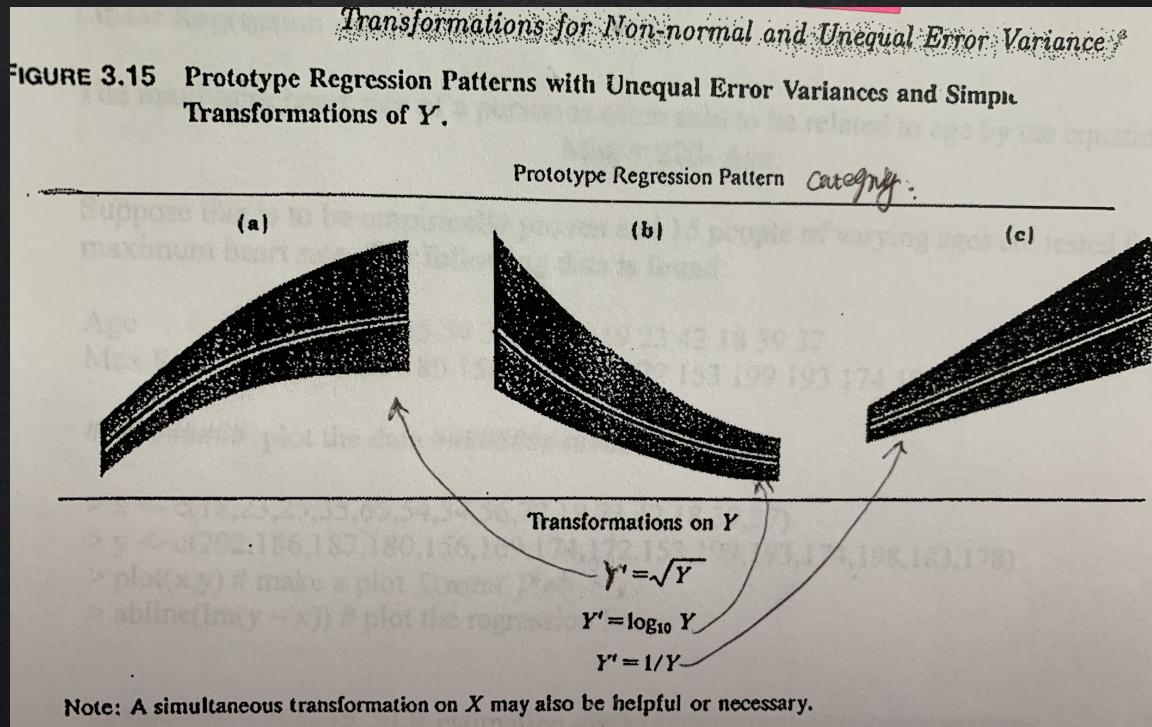
D S
X L

INDEPENDENT VARIABLE TRANSFORMATION

Prototype Nonlinear Regression Patterns with Constant Error Variance and Simple Transformations of X .



DEPENDENT VARIABLE TRANSFORMATION





MODEL PERFORMANCE EVALUATION

MSE: Mean square error

R-square(R2): R-squared is a goodness-of-fit measure for linear regression models. This statistic indicates the percentage of the variance in the dependent variable that the independent variables explain collectively. R-squared measures the strength of the relationship between your model and the dependent variable on a convenient 0 – 100% scale)

决定系数（英语： coefficient of determination，记为 R^2 或 r^2 ）在统计学中用于度量因变量的变异中可由自变量解释部分所占的比例，以此来判断统计模型的解释力。[\[1\]](#)[\[2\]](#)[\[3\]](#)

对于简单线性回归而言，决定系数为样本相关系数的平方。[\[4\]](#)当加入其他回归自变量后，决定系数相应地变为多重相关系数的平方。

假设一数据集包括 y_1, \dots, y_n 共 n 个观察值，相对应的模型预测值分别为 f_1, \dots, f_n 。定义残差 $e_i = y_i - f_i$ ，平均观察值为

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i.$$

于是可以得到总平方和

$$SS_{\text{tot}} = \sum_i (y_i - \bar{y})^2,$$

回归平方和

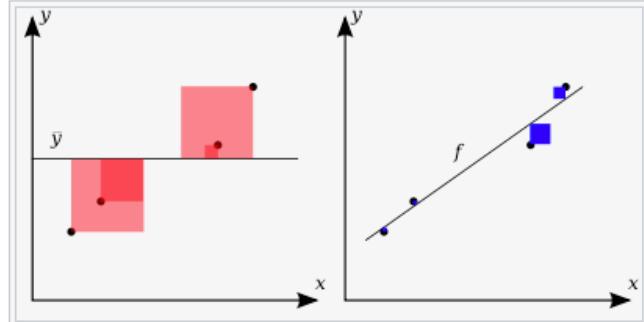
$$SS_{\text{reg}} = \sum_i (f_i - \bar{y})^2,$$

残差平方和

$$SS_{\text{res}} = \sum_i (y_i - f_i)^2 = \sum_i e_i^2$$

由此，决定系数可定义为

$$R^2 \equiv 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}.$$



决定系数 $R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$ 示意图 线性回归（右侧）的效果比起平均值（左侧）越好，决定系数的值就越接近于1。蓝色正方形表示线性回归的残差的平方，红色正方形数据表示对于平均值的残差的平方。



MAIN TAKEAWAY

1. “THE REGRESSION EQUATION MODELS THE RELATIONSHIP BETWEEN A RESPONSE VARIABLE Y AND A PREDICTOR VARIABLE X AS A LINE.”
2. “A REGRESSION MODEL YIELDS FITTED VALUES AND RESIDUALS PREDICTIONS OF THE RESPONSE AND THE ERRORS OF THE PREDICTIONS.”
3. “REGRESSION MODELS ARE TYPICALLY FIT BY THE METHOD OF LEAST SQUARES.”
4. “REGRESSION IS USED BOTH FOR PREDICTION AND EXPLANATION.”
5. A GOOD MODEL(ROBUST/APPLICATION RELEVANT) AND A REASONABLE EXPLANATION NEED FULLY MODEL INSPECTION



2.

LINEAR ALGEBRA 101

QUICK REVIEW OF BASIC LINEAR ALGEBRA

TWO DOCS HERE:

[Doc1](#)

[Doc2](#)



3.

MULTIPLE-REGRESSION

"THIS IS A SIMPLE EXPANSION OF SIMPLE LINEAR REGRESSION FROM 1 PREDICTOR TO MULTIPLE PREDICTORS"

$$Y = \beta_0 + \beta_1 X + \beta_2 X_2 + \dots + \beta_p X_p + \epsilon$$

$$\hat{Y} = b_0 + b_1 X_1 + b_2 X_2 + \dots + b_p X_p$$



MODEL ASSUMPTION AND MATRIX PRESENTATION

Multiple linear regression assumption:

1. Error is zero on average
2. Error has constant variance
3. Each observation is independent to other
4. ε_i are normal distributed

$$\begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} 1 & X_{2,1} & X_{3,1} \\ 1 & X_{2,2} & X_{3,2} \\ \vdots & \ddots & \ddots \\ 1 & X_{2,n} & X_{3,n} \end{bmatrix} \times \begin{bmatrix} B_1 \\ B_2 \\ B_3 \end{bmatrix} + \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix}$$
$$Y = X \times B + e$$



LEAST SQUARE ESTIMATION

$$Q(\underline{b}) = \sum e_i^2 = \sum_{i=1}^n (y_i - (b_0 + b_1 x_{i1} + b_2 x_{i2} + \cdots b_{p-1} x_{i,p-1}))^2$$

$$\sum e_i^2 = \mathbf{e}' \mathbf{e} = \mathbf{Y}' \mathbf{Y} - \mathbf{B}' \mathbf{X}' \mathbf{Y}$$

SOLUTION: FIND THE B THAT MINIMIZE Q(B)

$$\mathbf{B} = (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{Y} = \begin{bmatrix} B_1 \\ B_2 \\ B_3 \end{bmatrix}$$

$$\mathbf{X}' = \begin{bmatrix} 1 & 1 & \dots & 1 \\ X_{2,1} & X_{2,2} & \dots & X_{2,n} \\ X_{3,1} & X_{3,2} & \dots & X_{3,n} \end{bmatrix} \quad \mathbf{Y}' = \begin{bmatrix} Y_1 & Y_2 & \dots & Y_n \end{bmatrix} \quad \mathbf{B}' = \begin{bmatrix} B_1 & B_2 & B_3 \end{bmatrix}$$

$$\hat{\sigma}^2 = \frac{\mathbf{e}' \mathbf{e}}{n - k} = \text{residual variance}$$



COST FUNCTION

$$S(\beta_1, \beta_2, \beta_3) = \frac{1}{2N} \sum_{i=1}^n (y_i - (\beta_1 + \beta_2 x_2 + \beta_3 x_3))^2$$

GRADIENT DESCENT/PARTIAL DERIVATIVES

$$S'(\beta_1) = - (y - (\beta_1 + \beta_2 x_2 + \beta_3 x_3))$$

$$S'(\beta_2) = -x_2(y - (\beta_1 + \beta_2 x_2 + \beta_3 x_3))$$

$$S'(\beta_3) = -x_3(y - (\beta_1 + \beta_2 x_2 + \beta_3 x_3))$$



ACCESSING MODEL

Every Machine Learning method/Algorithm should have at least one objective function to optimize upon. And the metrics tied to this objective function will be the key performance indicator(KPI),

For linear model the most important KPI RMSE(Root Mean Square Error), which is the square root of the average squared error. This measures the overall accuracy of the model and is a basis for model comparison/selection

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}}$$

Other metrics: R², RSE, F-statistics, AIC, BIC



CROSS-VALIDATION

“Data Scientist/statistician beats their data to death to get what they want”

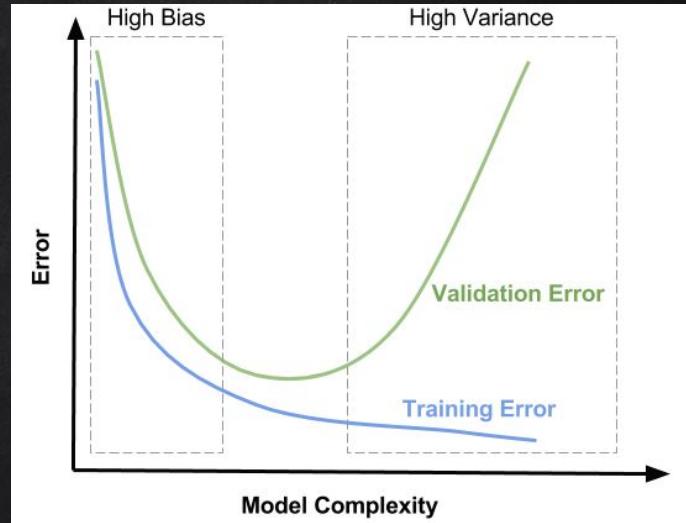
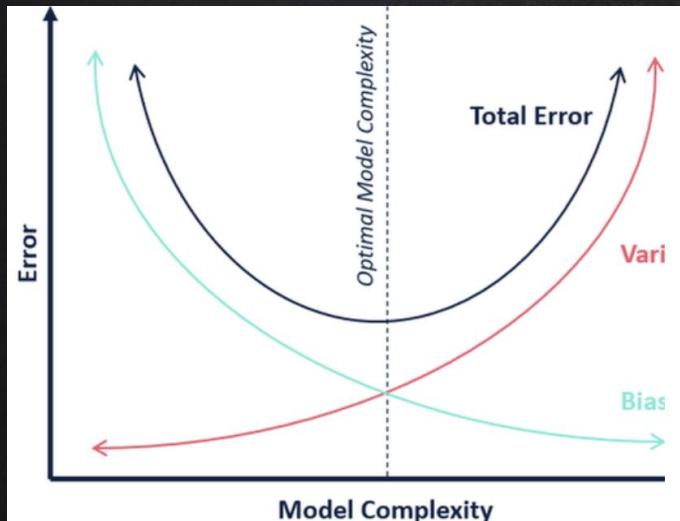
From certain point of view, it is a true statement. If a data was excessively used to fit a result, it is very easy to encounter another problem – Overfitting

What's the solution: holdout sample (build, test, validation), n-fold cross validation

Cross-validation extends the idea of a holdout sample to multiple sequential holdout samples. The algorithm for basic k-fold cross-validation is as follows:

1. Set aside $1/k$ of the data as a holdout sample.
2. Train the model on the remaining data.
3. Apply (score) the model to the $1/k$ holdout, and record needed model assessment
4. Restore the first $1/k$ of the data, and set aside the next $1/k$ (excluding any records that got picked the first time).
5. Repeat steps 2 and 3.
6. Repeat until each record has been used in the holdout portion.
7. Average or otherwise combine the model assessment metrics.

BIAS-VARIANCE TRADE-OFF



- 1) High Bias - underfitting
- 2) Goldilocks Zone - just right
- 3) High Variance - overfitting



MODEL SELECTION

A data may have many independent variables could be used for modeling, you may build the model with a handful predictors, but not necessary means the more variable you add to your model the better model you will get. If the model performance are similar, a simpler model is more preferred.

More predictor the smaller RMSE or high R-square, therefore adjusted R², AIC, BIC, F-statistics, used to select the best model with minimum predictors

$$R_{adj}^2 = 1 - (1 - R^2) \frac{n - 1}{n - P - 1}$$

Algorithm used to minimize AIC, BIC, maximize adjusted R² include backward selection, stepwise regression



WEIGHTED REGRESSION

Weighted regression used for variety of purposes;
in particular, it is important for analysis of complex surveys.

Data scientists may find weighted regression useful in two cases:

1. Inverse-variance weighting when different observations have been measured with different precision; the higher variance ones receiving lower weights.
2. Analysis of data where rows represent multiple cases; the weight variable encodes how many original observations each row represents.



KEY TAKEAWAY

1. “Multiple linear regression models the relationship between a response variable Y and multiple predictor variables ”.
2. “The most important metrics to evaluate a model are root mean squared error (RMSE/MSE) and R-squared (R²).”
3. “The standard error of the coefficients can be used to measure the reliability of a variable’s contribution to a model.”
4. “Stepwise regression is a way to automatically determine which variables should be included in the model.”
5. “Weighted regression is used to give certain records more or less weight in fitting the equation.”



LOGISTIC REGRESSION

"THIS IS ANALOGOUS ALGORITHM TO MLR WITH THE TARGET IS BINARY "

FOR EXAMPLE: IF X INCREASE THE ODDS OF Y=1 COMPARE TO Y=0 WILL INCREASE BY A CERTAIN RATIO OR DECREASE TO A CERTAIN RATIO



LOGISTIC FUNCTION

$$\rho = \beta_0 + \beta_1 X + \beta_2 X^2 + \dots + \beta_p X^p$$



$$\rho = 1 / 1 + \exp(-(\beta_0 + \beta_1 X + \beta_2 X^2 + \dots + \beta_p X^p))$$

Odds($Y=1$) = $\rho / (1 - \rho)$, therefore $\rho = \text{Odds} / (1 + \text{Odds})$

$$\text{Odds}($Y=1$) = \exp(\beta_0 + \beta_1 X + \beta_2 X^2 + \dots + \beta_p X^p)$$

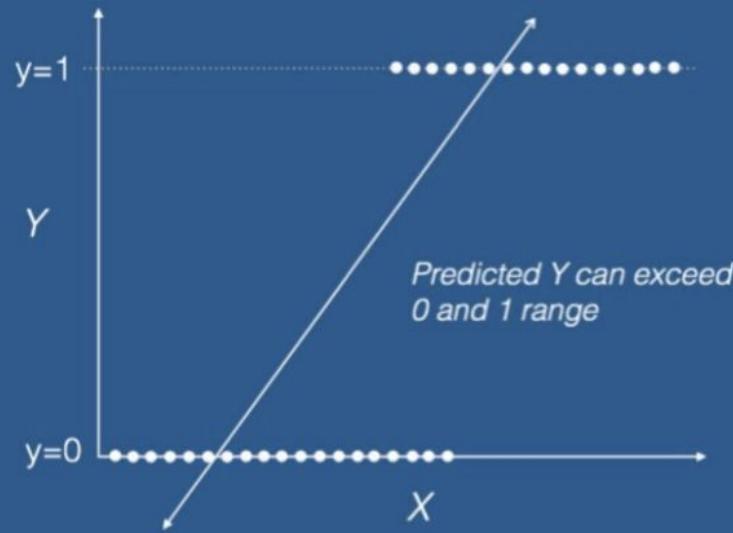
This is the linear function:

$$\log(\text{Odds}($Y=1$)) = \beta_0 + \beta_1 X + \beta_2 X^2 + \dots + \beta_p X^p$$

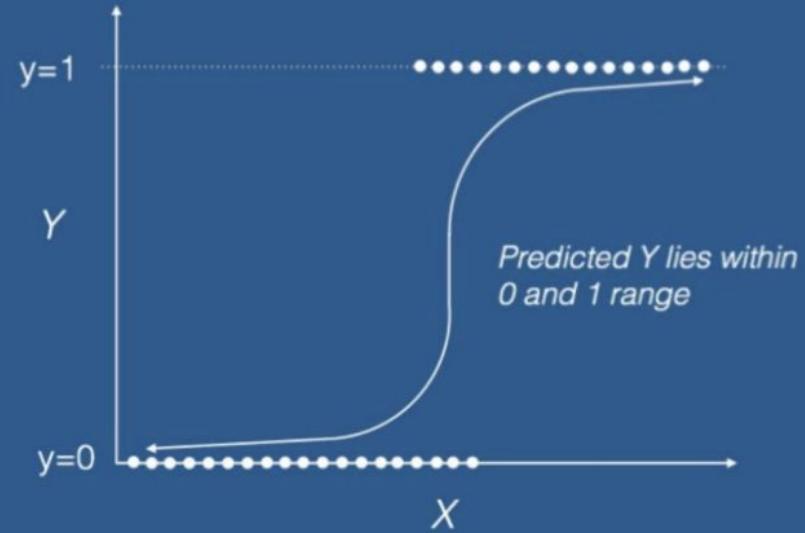
LINEAR REGRESSION VS LOGISTIC REGRESSION



Linear Regression



Logistic Regression





COST FUNCTION

$$J(\theta) = -\frac{1}{m} \sum \left[y^{(i)} \log(h\theta(x(i))) + (1 - y^{(i)}) \log(1 - h\theta(x(i))) \right]$$

Here,

$$h\theta(x) = h\beta(x) = \rho = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p)}}$$

Sigmoid function:

$$f(x) = \frac{1}{1 + e^{-(x)}}$$

MoreReading



EVALUATION OF CLASSIFICATION MODEL

Confusion matrix: "A tabular display (2×2 in the binary case) of the record counts by their predicted and actual classification status".

Sensitivity: "The percent (or proportion) of all 1s that are correctly classified as 1s. Also called Recall "

Specificity: "The percent (or proportion) of all 0s that are correctly classified as 0s."

Precision: "The percent (proportion) of predicted 1s that are actually 1s."

ROC curve: "A plot of sensitivity versus specificity."

AUC: "Area under a ROC curve"

Lift: "A measure of how effective the model is at identifying (comparatively rare) 1s at different probability cutoffs."

Accuracy: "The percent (or proportion) of cases classified correctly."

Other popular classification model evaluation metrics:
F1, KS, SomersD, GINI coefficient

CONFUSION MATRIX AND KEY CONCEPTS

Confusion matrix: “A tabular display (2×2 in the binary case) of the record counts by their predicted and actual classification status”.

		Predicted class	
		P	N
Actual Class	P	True Positives (TP)	False Negatives (FN)
	N	False Positives (FP)	True Negatives (TN)

Classification Evalu. Reading



5.

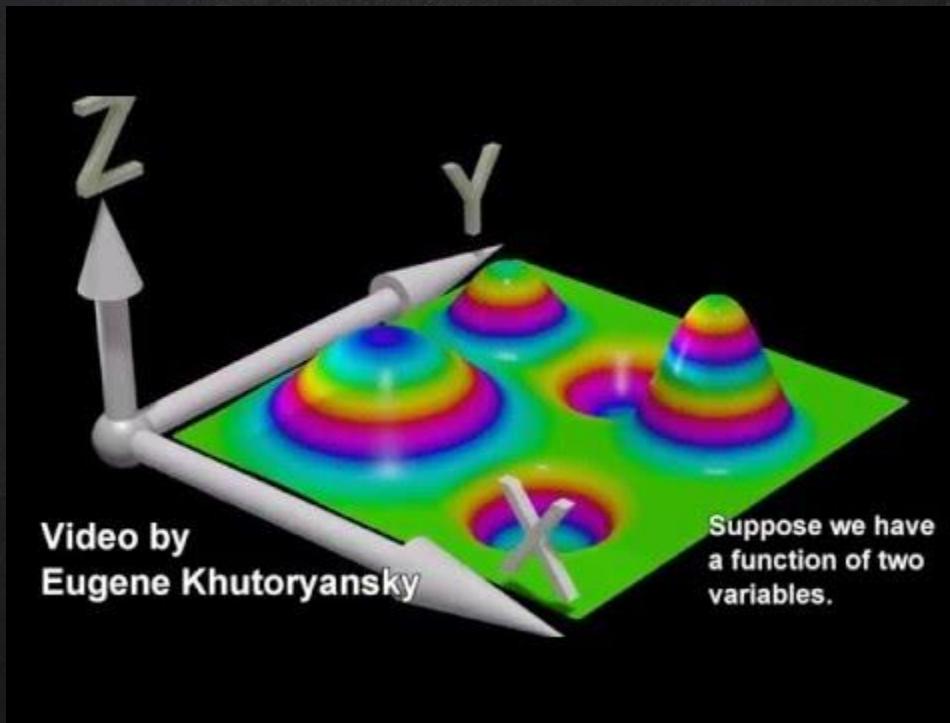
GRADIENT DESCENT

THIS ALGORITHM IS THE MOST POPULAR AND ONE FIT ALL METHOD USED IN MACHINE LEARNING TO SEARCHING FOR THE SOLUTION WHICH MINIMIZE THE COST FUNCTION. THIS IS A WAY TO FIND ANSWER OTHER THAN SOLVING A MATRIX EQUATION

THE KEY CONCEPTS ARE THE COST FUNCTION, (PARTIAL) DERIVATIVE, LEARNING RATE



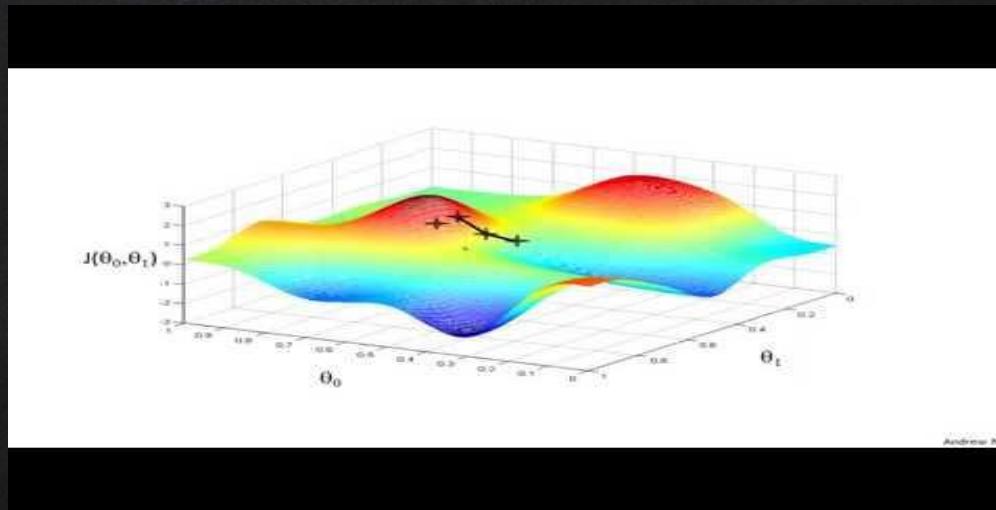
GRADIENT AND PARTIAL DERIVATIVE



From Physics Videos by Eugene Khutoryansky



GRADIENT DESCENT EXPLAINED/IMPLEMENTATION BY ANDREW NG



For whom have interested to learn more



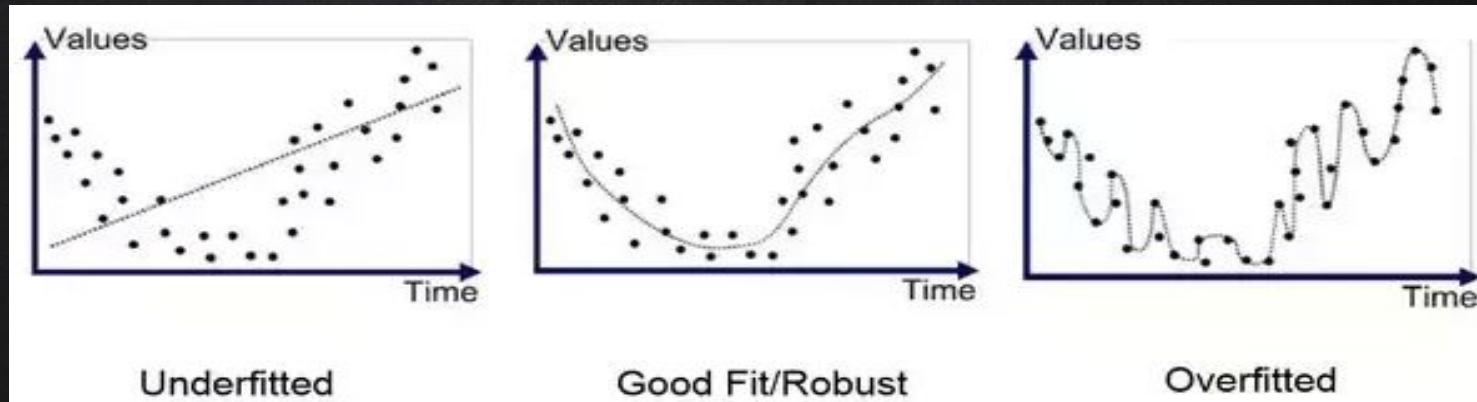
6.

REGULARIZATION

“IN MATHEMATICS, STATISTICS, AND COMPUTER SCIENCE, PARTICULARLY IN MACHINE LEARNING AND INVERSE PROBLEMS, REGULARIZATION IS THE PROCESS OF ADDING INFORMATION IN ORDER TO SOLVE AN ILL-POSED PROBLEM OR TO PREVENT OVERFITTING.[1]”

- WIKI

REGULARIZATION – OVERFITTING CONTROL



Three types of regularization in linear models, add regularization(penalized) term to cost function

Lasso regression (L1) : New Cost Function = Cost Function +

$$\lambda \sum_{j=1}^p |\beta_j|$$

Ridge regression (L2) : New Cost Function = Cost Function +

$$\lambda \sum_{j=1}^p \beta_j^2$$

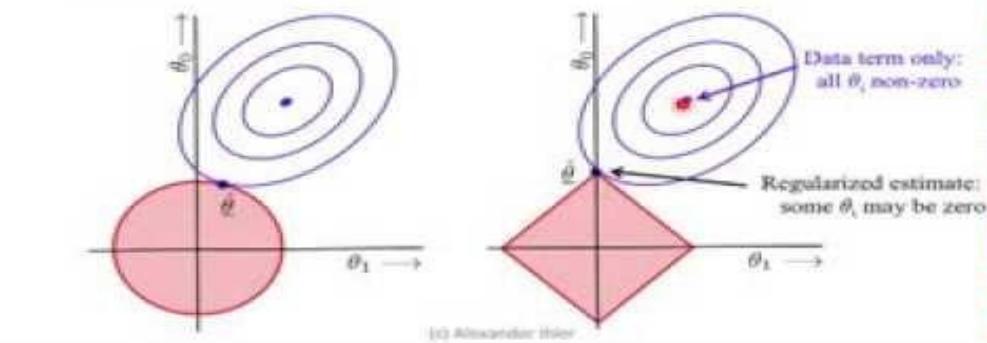
Elastic Net regression (L1+L2) : New Cost Function = Cost Function +

$$\lambda \sum_{j=1}^p |\beta_j| + \lambda \sum_{j=1}^p \beta_j^2$$

L1 AND L2 REGULARIZATION EXPLAINED

Regularization: L1 vs L2

- Estimate balances data term & regularization term
- Lasso tends to generate sparser solutions than a quadratic regularizer.



From Alexander Ihler channel



RECOMMENDED READINGS

<https://towardsdatascience.com/supervised-vs-unsupervised-learning-14f68e32ea8d>

<https://towardsdatascience.com/gradient-descent-algorithm-and-its-variants-10f652806a3>

<https://newbiettn.github.io/2016/08/30/precision-recall-sensitivity-specificity/>

<https://towardsdatascience.com/regularization-an-important-concept-in-machine-learning-5891628907ea>

<https://towardsdatascience.com/l1-and-l2-regularization-methods-ce25e7fc831c>

[Shared book and paper in github repo](#)



REFERENCES

<https://www.edureka.co/blog/statistics-and-probability/#Inferential%20Statistics>

<https://newbiettn.github.io/2016/08/30/precision-recall-sensitivity-specificity/>

<https://www.autodeskresearch.com/publications/samestats>

<https://www.edureka.co/blog/supervised-learning/#overview>

<https://towardsdatascience.com/introduction-to-machine-learning-algorithms-linear-regression-14c4e325882a>

<https://online.stat.psu.edu/stat462/node/93/>

[Bruce, Peter,Bruce, Andrew,Gedeck, Peter. Practical Statistics for Data Scientists \(Kindle Location 5267\). O'Reilly Media. Kindle Edition.](#)

[Physics Videos by Eugene Khutoryansky](#)

[Stanford university ,Machine Learning by Andrew Ng](#)

[https://en.wikipedia.org/wiki/Regularization_\(mathematics\)](https://en.wikipedia.org/wiki/Regularization_(mathematics))

<https://www.quora.com/What-is-regularization-in-machine-learning>