

Chapter7 Linear Models

From Risk Factors to Return Forecasts

Yanrui Zhou Wenchu Liu Shuxian Mei

Fintech Stu
Wuhan University

November10,2022



Linear Models

- 1 Multiple linear regression
- 2 How to build a linear factor model
- 3 Regularization and logistic regression



Multiple linear regression

The family of linear models represents one of the most useful hypothesis classes. Many learning algorithms that are widely applied in algorithmic trading rely on linear predictors.

- **efficiently trained**
- **robust to noisy financial data**
- **easy to interpret**
- **provide a good baseline**



How to formulate the linear model

In the population, the linear regression model has the following form for a single instance of the output y an input vector $\mathbf{X}^T = [x_1, \dots, x_p]$ and the error term ϵ :

$$y = f(\mathbf{x}) + \epsilon = \beta_0 + \beta_1 x_1 \dots + \beta_p x_p + \epsilon = \beta_0 + \sum_{j=1}^p \beta_j x_j + \epsilon$$

compactly, we have the matrix form:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$



Ordinary least squares

To minimize the sum of the squared distances between the output value and the value generated by the linear model.

$$\begin{aligned} \text{RSS}(\beta) &= \sum_{i=1}^N \epsilon_i^2 = \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 \\ &= (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) \end{aligned}$$

Thus, the least-squares coefficients β^{LS} are computed as:

$$\underset{\beta^{LS}}{\operatorname{argmin}} \text{RSS}(\beta) = (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta)$$



Ordinary least squares

$$\begin{aligned} \mathbf{y} &= N \times 1 & (\mathbf{y}^\top - \beta^\top \mathbf{X}^\top) (\mathbf{y} - \mathbf{X}\beta) \\ \mathbf{X} &= N \times p & = \mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X}\beta - \beta^\top \mathbf{X}^\top \mathbf{y} + \beta^\top \mathbf{X}^\top \mathbf{X}\beta \\ \beta &= p \times 1 & = \mathbf{y}^\top \mathbf{y} - 2\mathbf{y}^\top \mathbf{X}\beta + \beta^\top \mathbf{X}^\top \mathbf{X}\beta \end{aligned}$$

From two conclusions of matrix derivation:

$$\frac{\partial \mathbf{x}^\top \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^\top) \mathbf{x} \quad \frac{\partial \mathbf{A}^\top \mathbf{x}}{\partial \mathbf{x}} = \mathbf{A}$$

After derivation, we get:

$$\begin{aligned} -2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X}\beta &= 0 \\ \hat{\beta} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \end{aligned}$$



Ordinary least squares

$$\mathbf{y} = \mathbf{X}\beta + \epsilon = \hat{\mathbf{y}} + \epsilon$$

Geometric interpretation:

the coefficients that minimise RSS ensure that the vector of residuals $\mathbf{y} - \hat{\mathbf{y}}$ is orthogonal to the subspace of R^P spanned by the P columns of \mathbf{X} , and the estimates $\hat{\mathbf{y}}$ are orthogonal projections into that subspace.

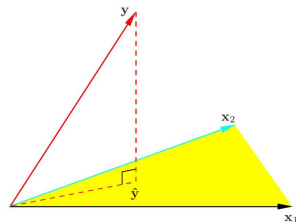


Figure: $\mathbf{y} = (\mathbf{x}_1, \mathbf{x}_2)\beta + \epsilon$



Ordinary least squares

prove : $\|y - \hat{y}\| < \|y - v\|$

$y - \hat{y}$ is orthogonal to plane W :

$$(y - \hat{y})(\hat{y} - v) = 0$$

$$\because y - v = (y - \hat{y}) + (\hat{y} - v)$$

$$\therefore \|y - v\|^2 = \|y - \hat{y}\|^2 + \|\hat{y} - v\|^2$$

$$\because \|\hat{y} - v\| \neq 0$$

$$\therefore \|y - v\|^2 - \|y - \hat{y}\|^2 > 0$$

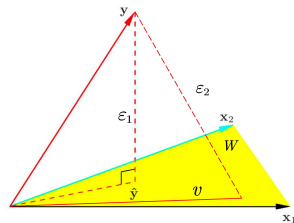


Figure: OLS prove



Maximum likelihood estimation

likelihood function: how likely it is to observe the sample of outputs when given the input as a function of the model parameters.

Let's set up the likelihood function by assuming a distribution for the error term, such as the standard normal distribution:

$$\epsilon_i \sim N(0, 1)$$

since normal distribution: $f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$

the conditional probability of observing a given output y_i :

$$p(y_i | x_i, \beta) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(y_i - x_i\beta)^2}{2\sigma^2}}$$



Maximum likelihood estimation

Assuming the output values are conditionally independent, given the inputs, we have:

$$L(y, x, \beta) = \prod_{i=1}^N p(y_i | x_i, \beta) = \left(\frac{1}{\sigma\sqrt{2\pi}} \right)^n e^{-\frac{\sum_{i=1}^N (y_i - x_i\beta)^2}{2\sigma^2}}$$

$$\log L(y, x, \beta) = n \log \left(\frac{1}{\sigma\sqrt{2\pi}} \right) - \frac{\sum_{i=1}^N (y_i - x_i\beta)^2}{2\sigma^2}$$



Gradient descent

Gradient descent is a general-purpose optimization algorithm that will find stationary points of smooth functions.gradient?

$$\left. \frac{\partial f}{\partial l} \right|_{(x_0, y_0)} = \lim_{\Delta s \rightarrow 0} \frac{f(x_0 + \Delta s \cos \alpha, y_0 + \Delta s \cos \beta) - f(x_0, y_0)}{\Delta s}$$

$$\stackrel{\text{differentiable}}{=} \lim_{\Delta s \rightarrow 0} \frac{f'_{x_0} \Delta s \cos \alpha + f'_{y_0} \Delta s \cos \beta + o(\Delta s)}{\Delta s}$$

$$\therefore \nabla f(x_0, y_0) = f'_{x_0} i + f'_{y_0} j \quad (\text{definition of Gradient})$$

$$\therefore \left. \frac{\partial f}{\partial l} \right|_{(x_0, y_0)} = \nabla f(x_0, y_0) \cdot e_l = f'_{x_0} \cos \alpha + f'_{y_0} \cos \beta$$

$$= |\nabla f(x_0, y_0)| |e_l| \cos \theta$$

It turns out that the maximal change of the function value results from a step in the direction of the gradient itself.



Gradient descent

current position
opposite direction
next position
small step
direction of fastest increase

$$\theta^1 = \theta^0 - \alpha \nabla J(\theta) \text{ evaluated at } \theta^0$$

https://log.csdn.net/q_41900366

minimise the cost of a prediction error:

$$J(\Theta) = \frac{1}{2N} \sum_{i=1}^N \left(h_{\theta} \left(x^{(i)} \right) - y^{(i)} \right)^2$$

$$h_{\theta} \left(x^{(i)} \right) = \Theta_0 + \Theta_1 x_1^{(i)}$$

$$\nabla J(\Theta) = \left\langle \frac{\delta J}{\delta \Theta_0}, \frac{\delta J}{\delta \Theta_1} \right\rangle$$

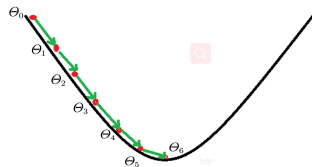


Figure: the process of GD



Gradient descent

For example, we have: $J(\Theta) = \theta_1^2 + \theta_2^2$ $\Theta^0 = (1, 3)$ $\alpha = 0.1$
then, $\nabla J(\Theta) = \langle 2\theta_1, 2\theta_2 \rangle$

$$\Theta^0 = (1, 3)$$

$$\Theta^1 = \Theta^0 - \alpha \nabla J(\Theta) = (1, 3) - 0.1 * (2, 6) = (0.8, 2.4)$$

$$\Theta^2 = (0.8, 2.4) - 0.1 * (1.6, 4.8) = (0.64, 1.92)$$

$$\Theta^3 = (0.5124, 1.536)$$

$$\Theta^4 = (0.4096, 1.2288000000000001)$$

$$\vdots$$

$$\Theta^{10} = (0.10737418240000003, 0.32212254720000005)$$

$$\Theta^{50} = (1.141798154164342e^{-05}, 3.42539442494306e^{-05})$$

$$\vdots$$

$$\Theta^{100} = (1.6296287810675902e^{-10}, 4.8888886343202771e^{-10})$$



The Gauss-Markov theorem

The baseline multiple regression has the following GMT assumptions:

- In the population, linearity holds so that $y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \epsilon$
- x_1, \dots, x_p satisfies I.I.D, then from the law of large numbers:

$$(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \epsilon = [\sum_{i=1}^N \mathbf{x}_i^\top \mathbf{x}_i]^{-1} [\sum_{i=1}^N \mathbf{x}_i^\top \epsilon_i] \xrightarrow{p} E[\mathbf{x}_i^\top \mathbf{x}_i]^{-1} E[\mathbf{x}_i^\top \epsilon_i] = 0$$
- No perfect collinearity : $r[\mathbf{X}^\top \mathbf{X}] = p$
- $E[\epsilon_i | \mathbf{x}_i] = 0$: $E[\hat{\beta} | \mathbf{X}] = \beta + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top E[\epsilon | \mathbf{X}] = \beta$
- Homoskedasticity : $E[\epsilon_i^2 | \mathbf{x}_i] = \sigma^2$: $\hat{\beta}$ is the **BLUE**

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{X}\beta + \epsilon) = \beta + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \epsilon$$

$$\mathbf{X} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{pmatrix}$$

$$\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{ip})$$

$$\mathbf{X}^\top = (\mathbf{x}_1^\top, \dots, \mathbf{x}_N^\top)$$

$$\mathbf{x}_i^\top \epsilon_i = (x_{i1}\epsilon_i, x_{i2}\epsilon_i, \dots, x_{ip}\epsilon_i)$$

$$\epsilon = \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_N \end{pmatrix}$$



How to conduct statistical inference

The key ingredient for statistical inference is a **test statistic with a known distribution** calculated from regression coefficients.

Adding the assumption of **normality**, we have:

$$\hat{\beta} \sim N(\beta, (\mathbf{X}^\top \mathbf{X})^{-1} \sigma^2) \quad (1)$$

$$\hat{\sigma}^2 = \frac{1}{N - p - 1} \sum_{i=1}^N (y_i - \hat{y}_i)^2 \quad (2)$$

$$t_j = \frac{\hat{\beta}_j}{\sqrt{v_j}} \sim t_{N-p-1} \quad (3)$$

$$(\mathbf{X}^\top \mathbf{X})^{-1} \sigma^2 = \begin{pmatrix} \text{var}(\beta_0) & \text{cov}(\beta_0, \beta_1) & \cdots \\ \text{cov}(\beta_1, \beta_0) & \ddots & \vdots \\ \vdots & \ddots & \vdots \\ \cdots & \cdots & \text{var}(\beta_{p+1}) \end{pmatrix}$$



How to diagnose and remedy problems

Diagnostics validate the model assumptions and help us prevent wrong conclusions when interpreting the result and conducting statistical inference. They include:

- Goodness of fit
- Remedy for Heteroskedasticity
- Serial correlation
- Multicollinearity



Goodness of fit

Goodness of fit metrics differ in how they measure the fit. For example:

- $R^2 = 1 - \frac{RSS}{TSS}$, and *adjusted* $R^2 = 1 - \frac{RSS/(N-p-1)}{TSS/(N-1)}$
- $AIC = -2 \log(\mathcal{L}^*) + 2p$
where \mathcal{L}^* is the value of the maximized likelihood function and p is the number of parameters.
- $BIC = -2 \log(\mathcal{L}^*) + \log(N)p$

Both metrics penalize for complexity. BIC imposes a higher penalty, so it might underfit relative to vice versa.

In practice, both criteria can be used jointly to guide model selection when the goal is an in-sample fit.



Heteroskedasticity

If the residual variance is correlated with an input variable:

- $\hat{\beta}$ is unbiased but is not the BLUE
- The statistical inference is not reliable any more.

Remedy for the heteroskedasticity problem under OLS estimates:

- **Robust standard errors** : $(X^T X)^{-1} X^T \Omega X (X^T X)^{-1}$
- **Clustered standard errors** : Assume there are distinct groups in the data that are homoscedastic, but the error variance differs between groups

We can also replace OLS by: WLS, GLSAR, GLS.....



Linear Models

- 1 Multiple linear regression
- 2 How to build a linear factor model
- 3 Regularization and logistic regression



The Basis of Multi-Factor Model

The average return of a stock is the payoff for taking risk.

- Average stock return = factor exposure \times factor premium:

$$r_i = \alpha_i + \beta_{i1}f_1 + \cdots + \beta_{iK}f_K + \epsilon_i$$

- And the average stock return is the cross-sum of the factor exposure and the factor premium:

$$E(r_i) = E(\alpha_i) + \beta_{i1}E(f_1) + \cdots + \beta_{iK}E(f_K)$$

- For the convenience of exposition, we use vectors to simplify the equation:

$$E(r_i) = \beta_i' E(f)$$



From the CAPM to the Fama-French factor models

- The CAPM model takes it that investors require compensation for systematic risk:

$$E[r_i] - r_f = \alpha_i + \beta_i (E[r_m] - r_f)$$

- The model implies that the value of α_i should be zero, which failed the empirical tests.
- Kenneth French and Eugene Fama identified additional risk factors:
 - **Size:** Market equity (ME)
 - **Value:** Book value of equity (BE) divided by ME
 - **Operating profitability (OP):** Revenue minus cost of goods sold/assets
 - **Investment:** Investment/assets



How Fama and French Obtain the factors

- The factors result from sorting stocks:

		BM		
		High(30%)	Middle(40%)	Low(30%)
ME	Small(50%)	S/H	S/M	S/L
	Big(50%)	B/H	B/M	B/L

- Then we get two factors:

$$SMB_t = \frac{(SL_t + SM_t + SH_t)}{3} - \frac{(BL_t + BM_t + BH_t)}{3}$$
$$HML_t = \frac{(BH_t + SH_t)}{2} - \frac{(BL_t + SL_t)}{2}$$

- Finally we run regression for risk exposure β_{HML} and β_{SMB}



Time Series Regression

- If risk factors are returns of a certain portfolio itself, we can use time series regression.
- Take CAPM model as an example:

ticker	time	$R_A - R_f$	$R_M - R_f$
A	2018	0.02	0.013
A	2019	-0.03	0.002
A	2020	0.02	0.02
A	2021	0.015	0.021

- We can run time series regressions to get the factor exposure:

$$R_{it} = \alpha_i + \beta'_i f_t + \varepsilon_{it}$$
$$E_T [R_i] = \alpha_i + \beta'_i E_T [f_t]$$



Cross-Section Regression

ticker	time	return	GDP(trillion)
A	2019	0.005	14
A	2020	0.003	14
A	2021	-0.002	17
B	2019	0.002	14
B	2020	-0.001	14
B	2021	-0.003	17
C	2019	0.012	14
C	2020	-0.002	14
C	2021	0.012	17

If risk factors are not returns themselves, we need one more step to get risk premium:

$$R_{it} = a_i + \beta'_i \lambda_t + \varepsilon_{it}$$
$$E_T [R_i] = \alpha_i + \beta'_i E_T [f_t]$$



Fama-Macbeth regression

ticker	time	return	GDP (trillion)	β
A	2019	0.005	14	β_A
A	2020	0.003	14	β_A
A	2021	-0.002	17	β_A
B	2019	0.002	14	β_B
B	2020	-0.001	14	β_B
B	2021	-0.003	17	β_B
C	2019	0.012	14	β_C
C	2020	-0.002	14	β_C
C	2021	0.012	17	β_C

First stage:

$$R_{it} = a_i + \beta_i' f_t + \varepsilon_{it}$$

Second stage:

$$R_{it} = \beta_i' \lambda_t + \alpha_{it}$$



How FM Regression Address the Inference Problem

If we use pool regression:

$$\begin{aligned}\hat{\beta} &= \frac{\sum_{i=1}^N \sum_{t=1}^T x_{it} y_{it}}{\sum_{i=1}^N \sum_{t=1}^T x_{it}^2} \\ &= \frac{\sum_{i=1}^N \sum_{t=1}^T x_{it} (x_{it} \beta + \varepsilon_{it})}{\sum_{i=1}^N \sum_{t=1}^T x_{it}^2} \\ &= \beta + \frac{\sum_{i=1}^N \sum_{t=1}^T x_{it} \varepsilon_{it}}{\sum_{i=1}^N \sum_{t=1}^T x_{it}^2}\end{aligned}$$



How FM Regression Address the Inference Problem

$$\begin{aligned}\text{Var}[\beta_{ols}] &= E \left(\left[\sum_{t=1}^T \sum_{i=1}^N x_{it} \varepsilon_{it} \right]^2 \left[\sum_{t=1}^T \sum_{i=1}^N x_{it}^2 \right]^{-2} \right) \\&= E \left(\left[\sum_{t=1}^T \left(\sum_{i=1}^N x_{it} \varepsilon_{it} \right)^2 \right] \left[\sum_{i=1}^N \sum_{t=1}^T x_{it}^2 \right]^{-2} \right) \\&= E \left(\left[\sum_{t=1}^T \left(\sum_{i=1}^N x_{it}^2 \varepsilon_{it}^2 + 2 \sum_{i=1}^N \sum_{s=i+1}^N x_{it} x_{st} \varepsilon_{it} \varepsilon_{st} \right) \right] \left[\sum_{i=1}^N \sum_{t=1}^T x_{it}^2 \right]^{-2} \right) \\&= (NT\sigma_x^2\sigma_\varepsilon^2 + TN(N-1)\rho_x\sigma_x^2\rho_\varepsilon\sigma_\varepsilon^2) / (NT\sigma_x^2)^2 \\&= \frac{\sigma_\varepsilon^2}{NT\sigma_x^2} (1 + (N-1)\rho_x\rho_\varepsilon)\end{aligned}$$



How FM Regression Address the Inference Problem

If we use cross-section regression:

$$\begin{aligned}\hat{\beta} &= \frac{\sum_{i=1}^N x_i E_T(y_i)}{\sum_{i=1}^N x_i^2} \\ &= \frac{\sum_{i=1}^N x_i (x_i \beta + \varepsilon_i)}{\sum_{i=1}^N x_i^2} \\ &= \beta + \frac{\sum_{i=1}^N x_i \varepsilon_i}{\sum_{i=1}^N x_i^2}\end{aligned}$$



How FM Regression Address the Inference Problem

$$\begin{aligned}\text{Var}[\beta_{ols}] &= E\left(\hat{\beta}_{ols} - \beta\right)^2 \\&= E\left(\left[\sum_{i=1}^N x_i \varepsilon_i\right]^2 \left[\sum_{i=1}^N x_i^2\right]^{-2}\right) \\&= E\left(\left[\sum_{i=1}^N x_i^2 \varepsilon_i^2 + \sum_{i=1}^{N-1} \sum_{j=i+1}^N x_i x_j \varepsilon_i \varepsilon_j\right] \left[\sum_{i=1}^N x_i^2\right]^{-2}\right) \\&= (N\sigma_x^2 \sigma_\varepsilon^2 + N(N-1)\rho_x \sigma_x^2 \rho_\varepsilon \sigma_\varepsilon^2) / (N\sigma_x^2)^2 \\&= \frac{\sigma_\varepsilon^2}{N\sigma_x^2} (1 + (N-1)\rho_x \rho_\varepsilon)\end{aligned}$$



How FM Regression Address the Inference Problem

If we use Fama-Macbeth regression:

$$\begin{aligned}\hat{\beta}_{FM} &= \sum_{t=1}^T \frac{\hat{\beta}_t}{T} \\ &= \frac{1}{T} \sum_{t=1}^T \left(\frac{\sum_{i=1}^N x_{it} y_{it}}{\sum_a^b x_{it}^2} \right) \\ &= \beta + \frac{1}{T} \sum_{t=1}^T \left(\frac{\sum_{i=1}^N x_{it} \varepsilon_{it}}{\sum_a^b x_{it}^2} \right)\end{aligned}$$



How FM Regression Address the Inference Problem

$$\begin{aligned}\text{Var}(\hat{\beta}_{FM}) &= \frac{1}{T^2} \text{Var}\left(\sum_{t=1}^T \hat{\beta}_t\right) \\&= \frac{\text{Var}(\hat{\beta}_t)}{T} + \frac{2 \sum_{t=1}^{T-1} \sum_{s=t}^T \text{Cov}(\hat{\beta}_t, \hat{\beta}_s)}{T^2} \\&= \frac{\text{Var}(\hat{\beta}_t)}{T} + \frac{T(T-1)}{T^2} \text{Cov}(\hat{\beta}_t, \hat{\beta}_s) \\&= \frac{1}{T} \left(\frac{\sigma_\epsilon^2}{N\sigma_X^2} \right) + \frac{T(T-1)}{T^2} * 0 \\&= \frac{\sigma_\epsilon^2}{NT\sigma_X^2}\end{aligned}$$



How FM Regression Address the Inference Problem

- The stand error from pool regression is:

$$\text{Var}[\beta_{ols}] = \frac{\sigma_{\epsilon}^2}{N\sigma_x^2} (1 + (N-1)\rho_x\rho_{\epsilon})$$

- The stand error from cross-section is:

$$\text{Var}[\beta_{ols}] = \frac{\sigma_{\epsilon}^2}{N\sigma_x^2} (1 + (N-1)\rho_x\rho_{\epsilon})$$

- The stand error from Fama-Macbeth is:

$$\text{Var}(\hat{\beta}_{FM}) = \frac{\sigma_{\epsilon}^2}{NT\sigma_x^2}$$



Linear Models

- 1 Multiple linear regression
- 2 How to build a linear factor model
- 3 Regularization and logistic regression



Regularizing linear regression using shrinkage

When a linear regression model contains many correlated variables, their coefficients will be poorly determined, and the model will more likely to overfit the sample.

- In OLS, β is estimated like this:

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- Notice that $(\mathbf{X}^T \mathbf{X})^{-1}$ equals $\frac{(\mathbf{X}^T \mathbf{X})^*}{|\mathbf{X}^T \mathbf{X}|}$



Shrinkage models

One popular technique to control overfitting is regularization, which involves the addition of a penalty term to the error function to discourage the coefficients from reaching large values.

$$\hat{\beta}^S = \underset{\beta^S}{\operatorname{argmin}} \sum_{i=1}^N \left[\left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_j \right)^2 + \lambda S(\beta) \right]$$



Ridge Regression

- The penalty term:

$$S(\beta) = \sum_{i=1}^p \beta_i^2 = \|\beta\|^2$$

- Hence, the ridge coefficients are defined as:

$$\begin{aligned}\hat{\beta}^{\text{Ridge}} &= \underset{\beta^{\text{Ridge}}}{\operatorname{argmin}} \sum_{i=1}^N \left[\left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_j \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right] \\ &= \underset{\beta^{\text{Ridge}}}{\operatorname{argmin}} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) + \lambda \beta^T \beta\end{aligned}$$



Ridge Regression

- The intercept β_0 has been excluded from the penalty to make the procedure independent of the origin chosen for the output variable.
- It is important to standardize the inputs. This is because the ridge solution is sensitive to the scale of the inputs.
- There is also a closed solution for the ridge estimator that resembles the OLS case:

$$\hat{\beta}^{\text{Ridge}} = \left(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I} \right)^{-1} \mathbf{X}^T \mathbf{y}$$



Lasso Regression

The lasso penalty is the sum of the absolute values of the coefficient vector.

- Hence, the lasso coefficients are defined as:

$$\begin{aligned}\hat{\beta}^{\text{Lasso}} &= \underset{\beta^{\text{Lasso}}}{\operatorname{argmin}} \sum_{i=1}^N \left[\left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p |\beta_j| \right] \\ &= \underset{\beta^{\text{Lasso}}}{\operatorname{argmin}} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) + \lambda \|\beta\|_1\end{aligned}$$



Lasso Regression

- Similar to ridge regression, the inputs need to be standardized. The lasso penalty makes the solution nonlinear.
- The lasso penalty had the effect of gradually reducing some coefficients to **zero** as the regularization increases. For this reason, the lasso can be used for the continuous selection of a subset of features.



Logistic Regression

- The output variable y :

$$y_t = \begin{cases} 1 & A \\ 0 & \text{otherwise} \end{cases}$$

- Logistic regression models the probability that y belongs to either of the categories.

$$P(\mathbf{x}_t) = P_r(y_t = 1 \mid \mathbf{x}_t).$$



Logistic Regression

To prevent the model from producing values outside the $[0, 1]$ interval, we must model $p(x)$ using a function that only gives outputs between 0 and 1 over the entire domain of x .

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



Logistic Regression

More formally, we are seeking to maximize the likelihood function:

$$\max_{\beta} \mathcal{L}(\beta) = \prod_{i: y_i=1} p(\mathbf{x}_i) \prod_{i': y_{i'}=0} (1 - p(\mathbf{x}_{i'}))$$

$$\beta^{\text{ML}} = \text{argmax} \log \mathcal{L}(\beta)$$

$$\text{i.e.} \sum_{i=1}^N (y_i \log p(\mathbf{x}_i, \beta) + (1 - y_i) \log (1 - p(\mathbf{x}_i, \beta)))$$



Confusion matrix

		Actual (Truth)				
		Positive	Negative	Accuracy	$= \frac{\text{\# Correct Predictions}}{\text{\# Cases}} = \frac{TP + TN}{TP + FP + TN + FN}$	
Prediction	Positive	True Positive (TP)	False Positive (FP)	True Positive Rate (Sensitivity, Recall)	$= \frac{\text{\# Correct Positive Predictions}}{\text{\# Positive Cases}} = \frac{TP}{TP + FN}$	
	Negative	False Negative (FN)	True Negative (TN)	False Negative Rate (Miss Rate)	$= 1 - \text{True Positive Rate}$	
				True Negative Rate (Specificity)	$= \frac{\text{\# Correct Negative Predictions}}{\text{\# Negative Cases}} = \frac{TN}{TN + FP}$	
				False Positive Rate (Fall-Out)	$= 1 - \text{True Negative Rate}$	

Figure: Confusion matrix and related error metrics



ROC and AUC

- Receiver operating characteristics (ROC) curve allows us to visualize and select classifiers based on their performance.
- The area under the curve (AUC) is defined as the area under the ROC plot that varies between 0.5 and the maximum of 1.

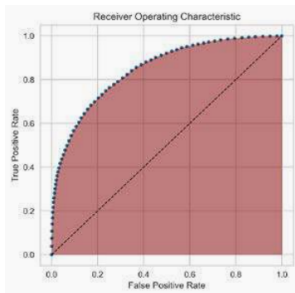


Figure: Receiver operating characteristics (ROC) curve

