

Regression

MLAI: Week 3

Neil D. Lawrence

Department of Computer Science
Sheffield University

13th October 2015

Review

- ▶ Last time: Looked at objective functions for movie recommendation.
- ▶ Minimized sum of squares objective by steepest descent and stochastic gradients.
- ▶ This time: explore least squares for regression.

Outline

Regression

Regression Examples

- ▶ Predict a real value, y_i given some inputs x_i .
- ▶ Predict quality of meat given spectral measurements (Tecator data).
- ▶ Radiocarbon dating, the C14 calibration curve: predict age given quantity of C14 isotope.
- ▶ Predict quality of different Go or Backgammon moves given expert rated training data.

Olympic 100m Data

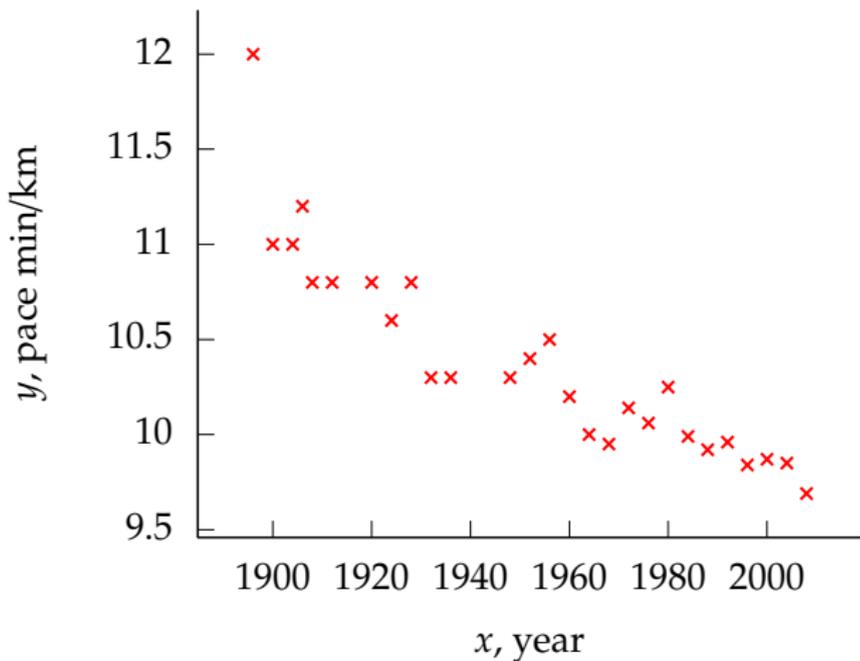
- ▶ Gold medal times for Olympic 100 m runners since 1896.



Image from Wikimedia
Commons

<http://bit.ly/191adDC>

Olympic 100m Data



Olympic 100 m Data.

Olympic Marathon Data

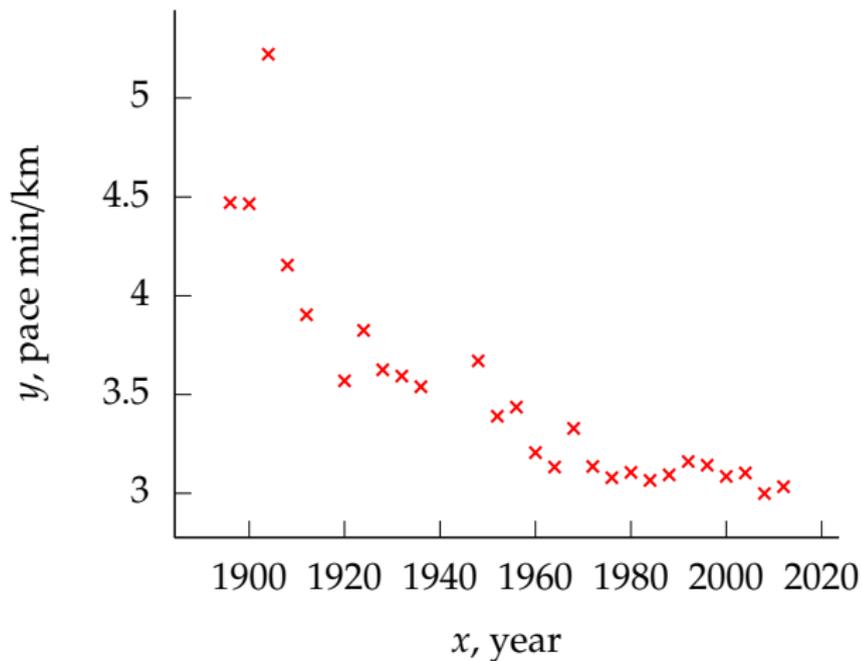
- ▶ Gold medal times for Olympic Marathon since 1896.
- ▶ Marathons before 1924 didn't have a standardised distance.
- ▶ Present results using pace per km.
- ▶ In 1904 Marathon was badly organised leading to very slow times.



Image from Wikimedia
Commons

<http://bit.ly/16kMKHQ>

Olympic Marathon Data



Olympic Marathon Data.

What is Machine Learning?

data

- ▶ **data**: observations, could be actively or passively acquired (meta-data).

What is Machine Learning?

data +

- ▶ **data**: observations, could be actively or passively acquired (meta-data).

What is Machine Learning?

data + **model**

- ▶ **data**: observations, could be actively or passively acquired (meta-data).
- ▶ **model**: assumptions, based on previous experience (other data! transfer learning etc), or beliefs about the regularities of the universe. Inductive bias.

What is Machine Learning?

data + **model** =

- ▶ **data**: observations, could be actively or passively acquired (meta-data).
- ▶ **model**: assumptions, based on previous experience (other data! transfer learning etc), or beliefs about the regularities of the universe. Inductive bias.

What is Machine Learning?

$$\text{data} + \text{model} = \text{prediction}$$

- ▶ **data**: observations, could be actively or passively acquired (meta-data).
- ▶ **model**: assumptions, based on previous experience (other data! transfer learning etc), or beliefs about the regularities of the universe. Inductive bias.
- ▶ **prediction**: an action to be taken or a categorization or a quality score.

Regression: Linear Relationship

$$y = mx + c$$

- ▶ y : winning time/pace.

Regression: Linear Relationship

$$y = mx + c$$

- ▶ y : winning time/pace.
- ▶ x : year of Olympics.

Regression: Linear Relationship

$$y = mx + c$$

- ▶ y : winning time/pace.
- ▶ x : year of Olympics.
- ▶ m : rate of improvement over time.

Regression: Linear Relationship

$$y = mx + c$$

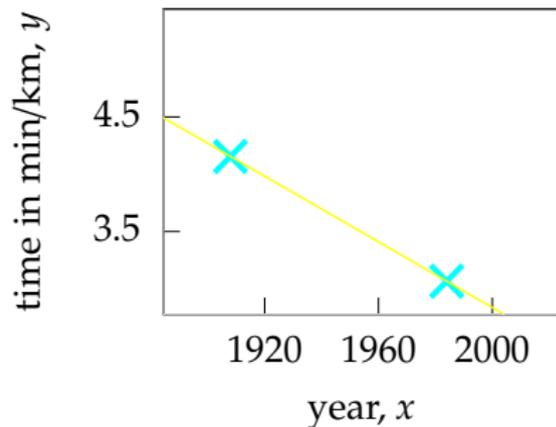
- ▶ y : winning time/pace.
- ▶ x : year of Olympics.
- ▶ m : rate of improvement over time.
- ▶ c : winning time at year 0.

Two Simultaneous Equations

A system of two simultaneous equations with two unknowns.

$$y_1 = mx_1 + c$$

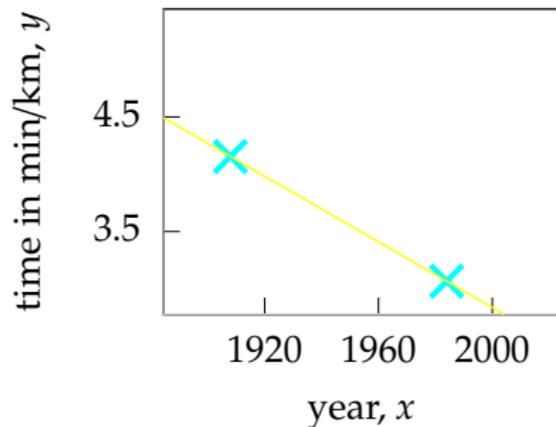
$$y_2 = mx_2 + c$$



Two Simultaneous Equations

A system of two simultaneous equations with two unknowns.

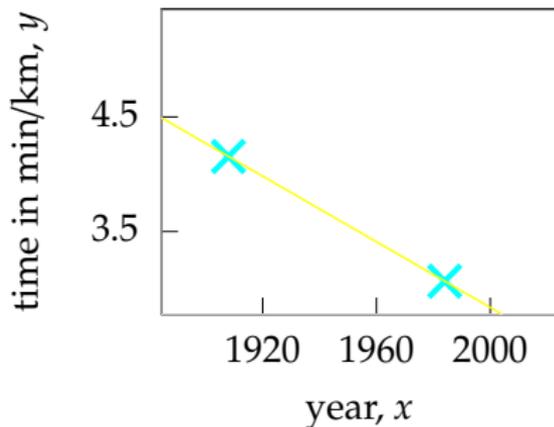
$$y_1 - y_2 = m(x_1 - x_2)$$



Two Simultaneous Equations

A system of two simultaneous equations with two unknowns.

$$\frac{y_1 - y_2}{x_1 - x_2} = m$$

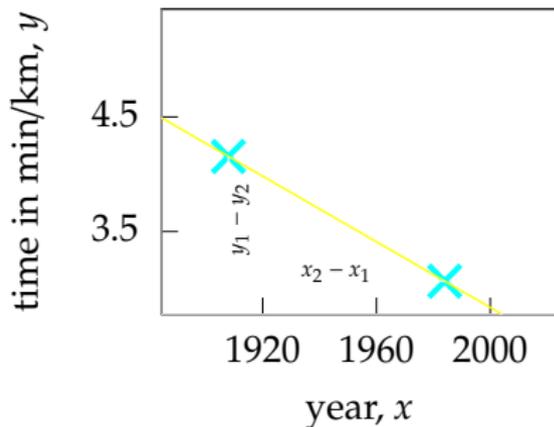


Two Simultaneous Equations

A system of two simultaneous equations with two unknowns.

$$m = \frac{y_2 - y_1}{x_2 - x_1}$$

$$c = y_1 - mx_1$$



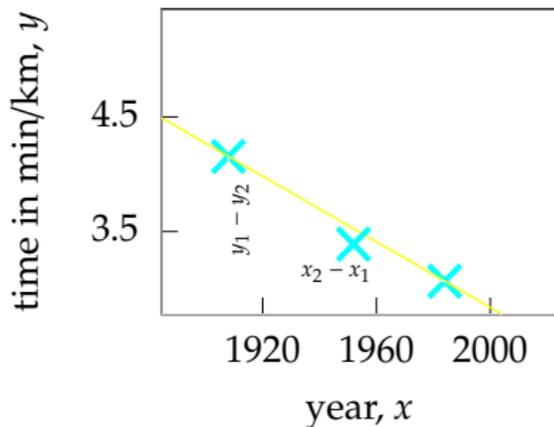
Two Simultaneous Equations

How do we deal with three simultaneous equations with only two unknowns?

$$y_1 = mx_1 + c$$

$$y_2 = mx_2 + c$$

$$y_3 = mx_3 + c$$



Overdetermined System

- ▶ With two unknowns and two observations:

$$y_1 = mx_1 + c$$

$$y_2 = mx_2 + c$$

Overdetermined System

- ▶ With two unknowns and two observations:

$$y_1 = mx_1 + c$$

$$y_2 = mx_2 + c$$

- ▶ Additional observation leads to *overdetermined* system.

$$y_3 = mx_3 + c$$

Overdetermined System

- ▶ With two unknowns and two observations:

$$y_1 = mx_1 + c$$

$$y_2 = mx_2 + c$$

- ▶ Additional observation leads to *overdetermined* system.

$$y_3 = mx_3 + c$$

- ▶ This problem is solved through a noise model $\epsilon \sim \mathcal{N}(0, \sigma^2)$

$$y_1 = mx_1 + c + \epsilon_1$$

$$y_2 = mx_2 + c + \epsilon_2$$

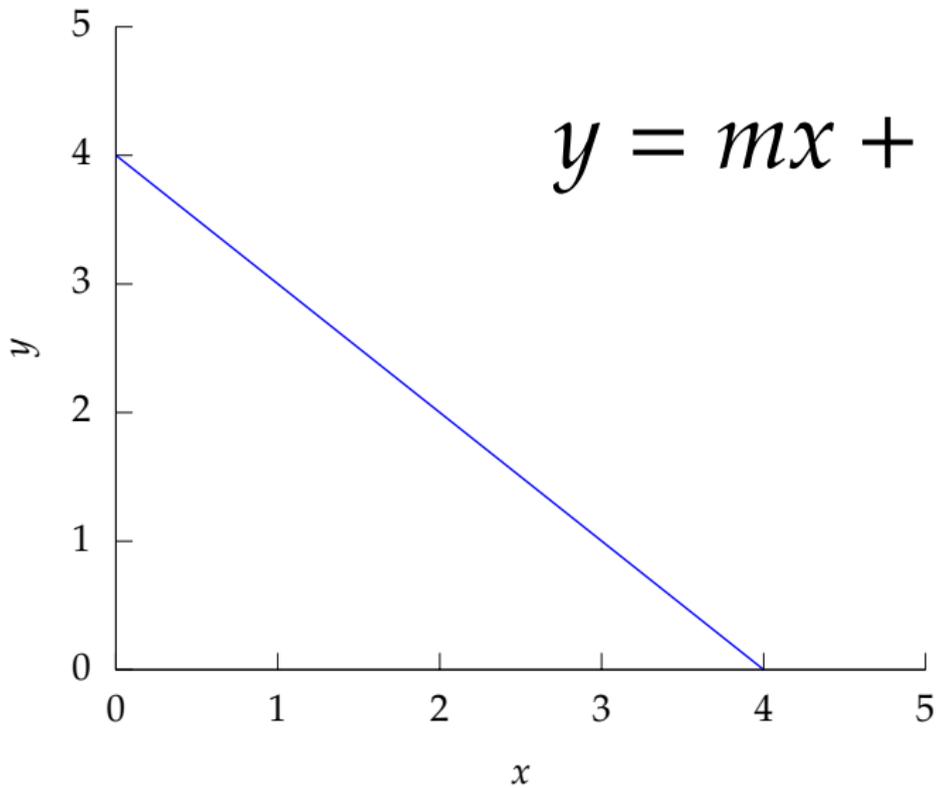
$$y_3 = mx_3 + c + \epsilon_3$$

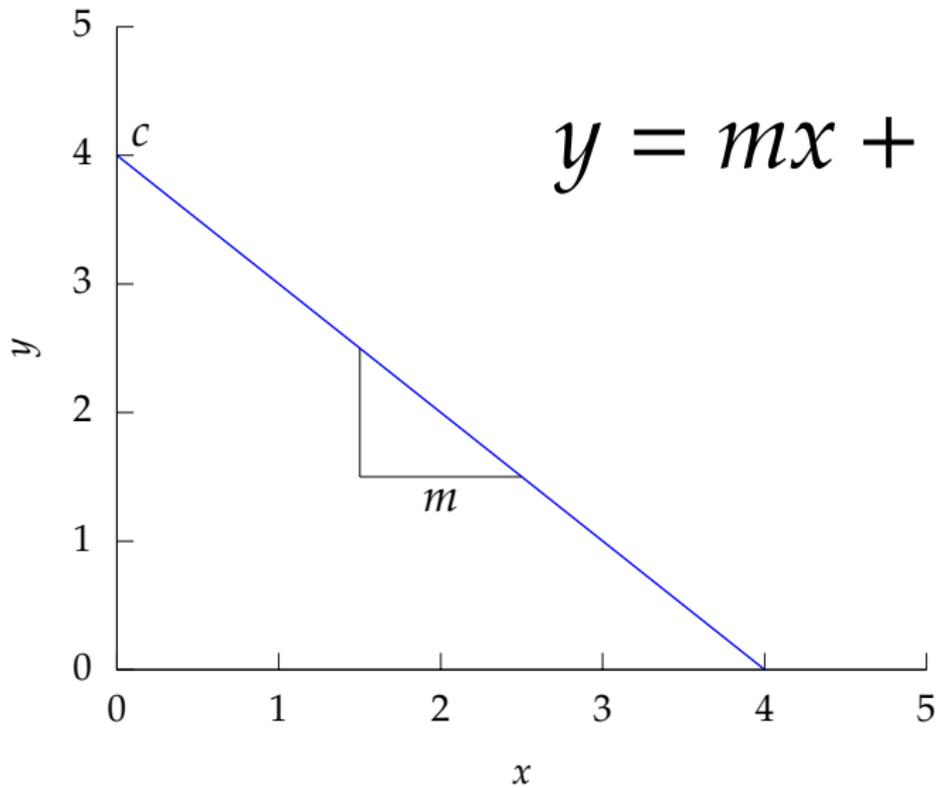
Noise Models

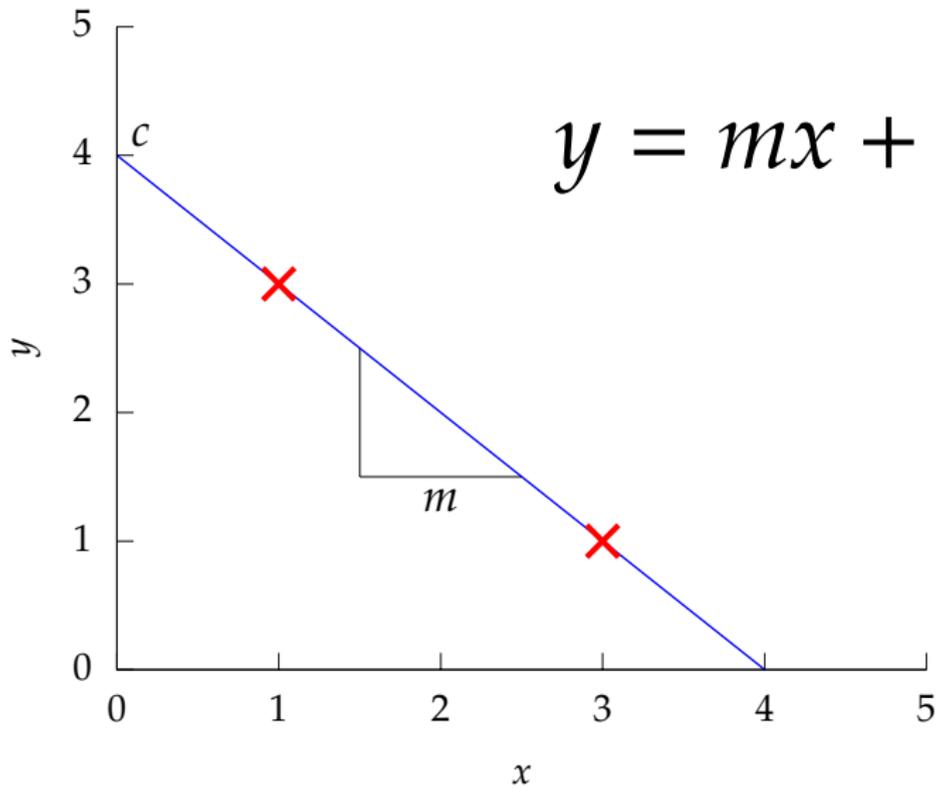
- ▶ We aren't modeling entire system.
- ▶ Noise model gives mismatch between model and data.
- ▶ Gaussian model justified by appeal to central limit theorem.
- ▶ Other models also possible (Student- t for heavy tails).
- ▶ Maximum likelihood with Gaussian noise leads to *least squares*.

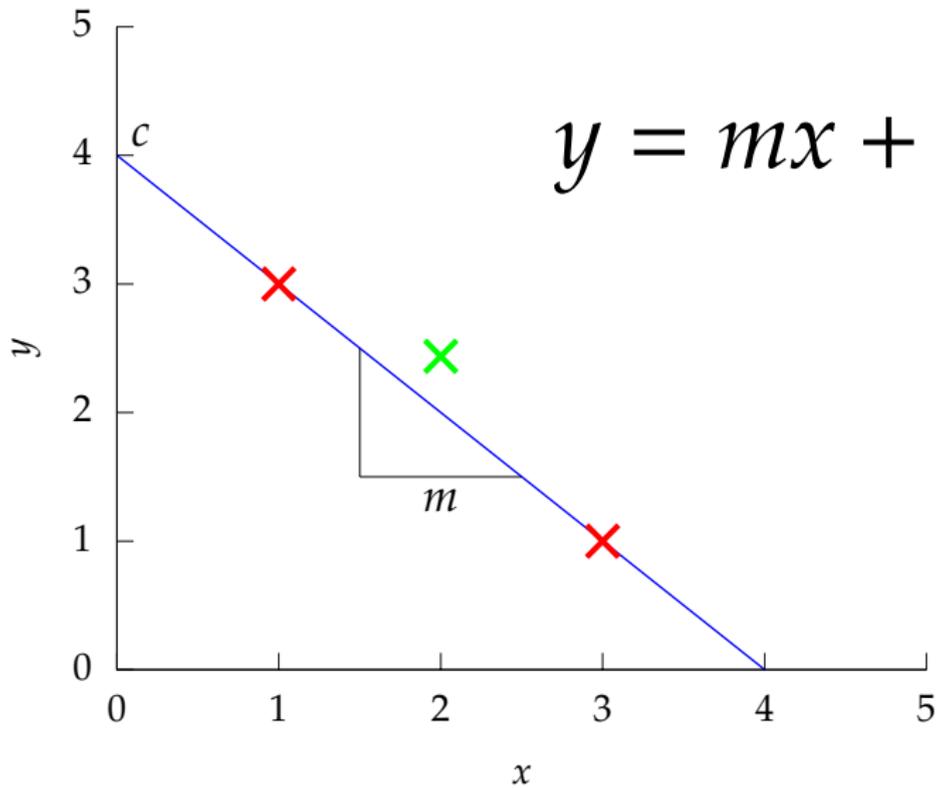
$$y = mx + c$$

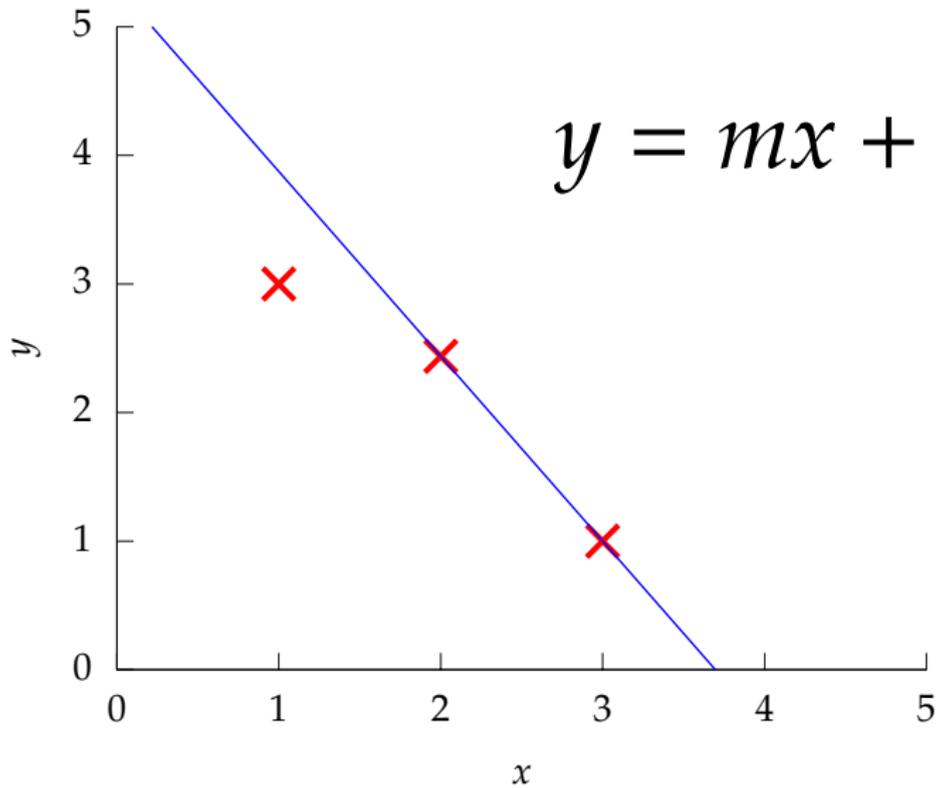
$$y = mx + c$$



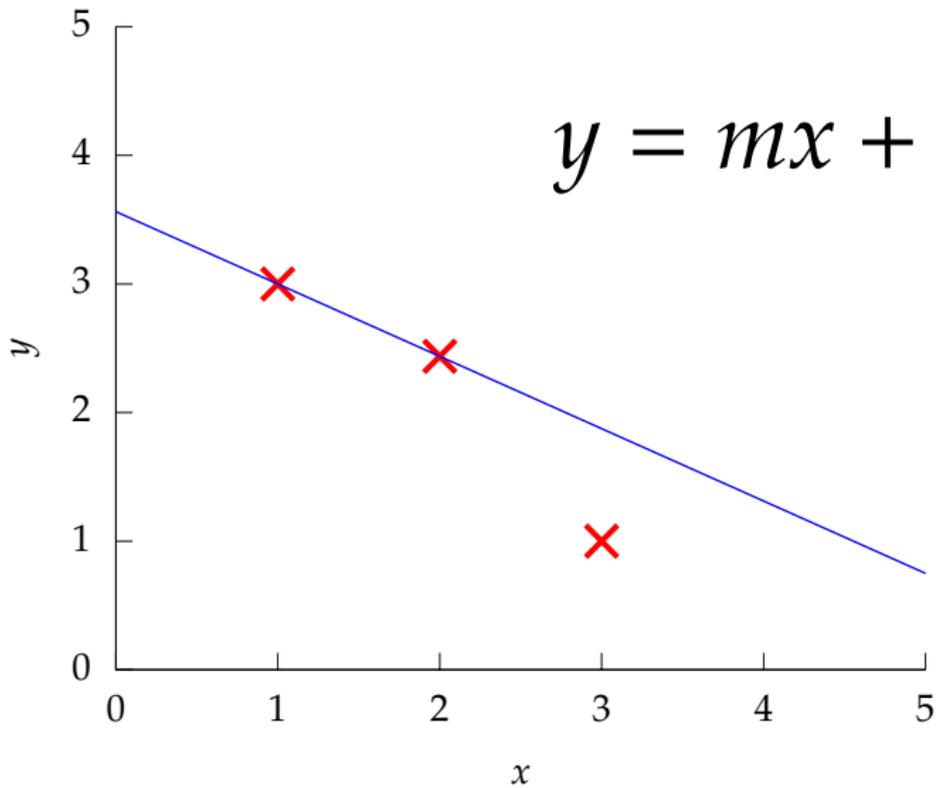


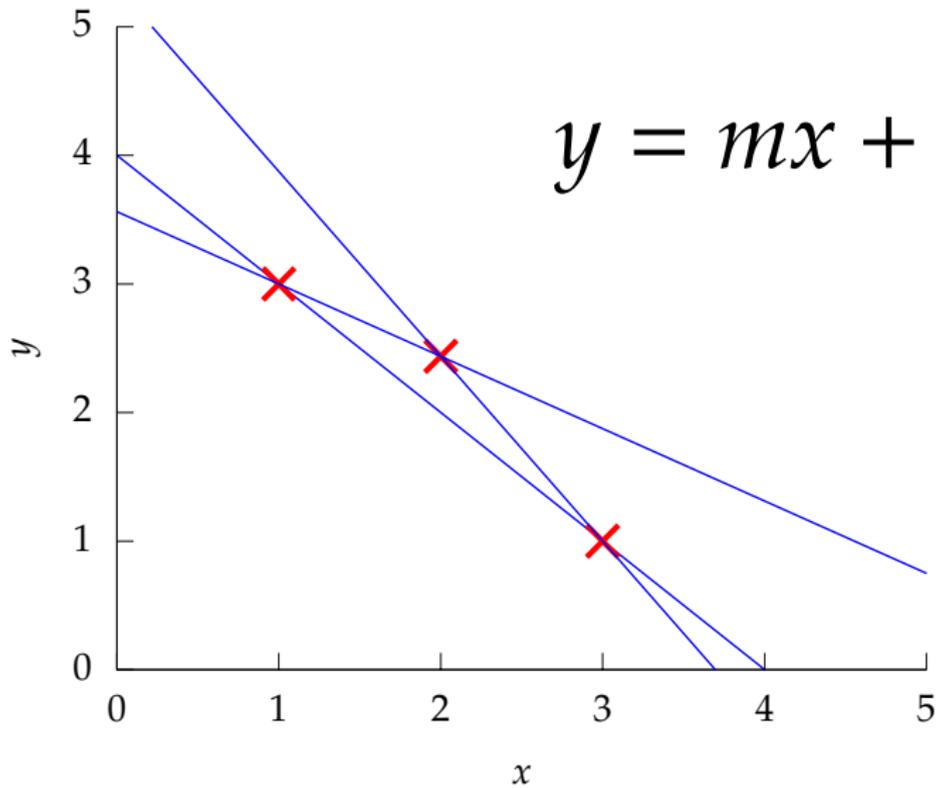






$$y = mx + c$$





$$y = mx + c$$

point 1: $x = 1, y = 3$

$$3 = m + c$$

point 2: $x = 3, y = 1$

$$1 = 3m + c$$

point 3: $x = 2, y = 2.5$

$$2.5 = 2m + c$$



riens. L'opinion contraire est une illusion de l'esprit qui, perdant de vue les raisons fugitives du choix de la volonté dans les choses indifférentes, se persuade qu'elle s'est déterminée d'elle-même et sans motifs.

Nous devons donc envisager l'état présent de l'univers, comme l'effet de son état antérieur, et comme la cause de celui qui va suivre. Une intelligence qui, pour un instant donné, connaîtrait toutes les forces dont la nature est animée, et la situation respective des êtres qui la composent, si d'ailleurs elle était assez vaste pour soumettre ces données à l'analyse, embrasserait dans la même formule les mouvemens des plus grands corps de l'univers et ceux du plus léger atome : rien ne serait incertain pour elle, et l'avenir comme le passé, serait présent à ses yeux. L'esprit humain offre, dans la perfection qu'il a su donner à l'Astronomie, une faible esquisse de cette intelligence. Ses découvertes en Mécanique et en Géométrie, jointes à celle de la pesanteur universelle, l'ont mis à portée de comprendre dans les mêmes expressions analytiques, les états passés et futurs du système du monde. En appliquant la même méthode à quelques autres objets de ses connaissances, il est parvenu à ramener à des lois générales, les phénomènes observés, et à prévoir ceux que des circonstances données doivent faire éclore. Tous ces efforts dans la recherche de la vérité, tendent à le rapprocher sans cesse de l'intelligence que nous venons de concevoir, mais dont il restera toujours infiniment éloigné. Cette tendance propre à l'espèce humaine, est ce qui la rend supérieure aux animaux; et ses progrès en ce genre, distinguent les nations et les siècles, et font leur véritable gloire.

Rappelons-nous qu'autrefois, et à une époque qui

other, we say that its choice is an effect without a cause. It is then, says Leibnitz, the blind chance of the Epicureans. The contrary opinion is an illusion of the mind, which, losing sight of the evasive reasons of the choice of the will in indifferent things, believes that choice is determined of itself and without motives.

We ought then to regard the present state of the universe as the effect of its anterior state and as the cause of the one which is to follow. Given for one instant an intelligence which could comprehend all the forces by which nature is animated and the respective situation of the beings who compose it—an intelligence sufficiently vast to submit these data to analysis—it would embrace in the same formula the movements of the greatest bodies of the universe and those of the lightest atom; for it, nothing would be uncertain and the future, as the past, would be present to its eyes. The human mind offers, in the perfection which it has been able to give to astronomy, a feeble idea of this intelligence. Its discoveries in mechanics and geometry, added to that of universal gravity, have enabled it to comprehend in the same analytical expressions the past and future states of the system of the world. Applying the same method to some other objects of its knowledge, it has succeeded in referring to general laws observed phenomena and in foreseeing those which given circumstances ought to produce. All these efforts in the search for truth tend to lead it back continually to the vast intelligence which we have just mentioned, but from which it will always remain infinitely removed. This tendency, peculiar to the human race, is that which renders it superior to animals; and their progress

height: "The day will come when, by study pursued through several ages, the things now concealed will appear with evidence; and posterity will be astonished that truths so clear had escaped us." Clairaut then undertook to submit to analysis the perturbations which the comet had experienced by the action of the two great planets, Jupiter and Saturn; after immense calculations he fixed its next passage at the perihelion toward the beginning of April, 1759, which was actually verified by observation. The regularity which astronomy shows us in the movements of the comets doubtless exists also in all phenomena. .

The curve described by a simple molecule of air or vapor is regulated in a manner just as certain as the planetary orbits; the only difference between them is that which comes from our ignorance.

Probability is relative, in part to this ignorance, in part to our knowledge. We know that of three or a greater number of events a single one ought to occur; but nothing induces us to believe that one of them will occur rather than the others. In this state of indecision it is impossible for us to announce their occurrence with certainty. It is, however, probable that one of these events, chosen at will, will not occur because we see several cases equally possible which exclude its occurrence, while only a single one favors it.

The theory of chance consists in reducing all the events of the same kind to a certain number of cases equally possible, that is to say, to such as we may be equally undecided about in regard to their existence, and in determining the number of cases favorable to the event whose probability is sought. The ratio of

$$y = mx + c + \epsilon$$

point 1: $x = 1, y = 3$

$$3 = m + c + \epsilon_1$$

point 2: $x = 3, y = 1$

$$1 = 3m + c + \epsilon_2$$

point 3: $x = 2, y = 2.5$

$$2.5 = 2m + c + \epsilon_3$$

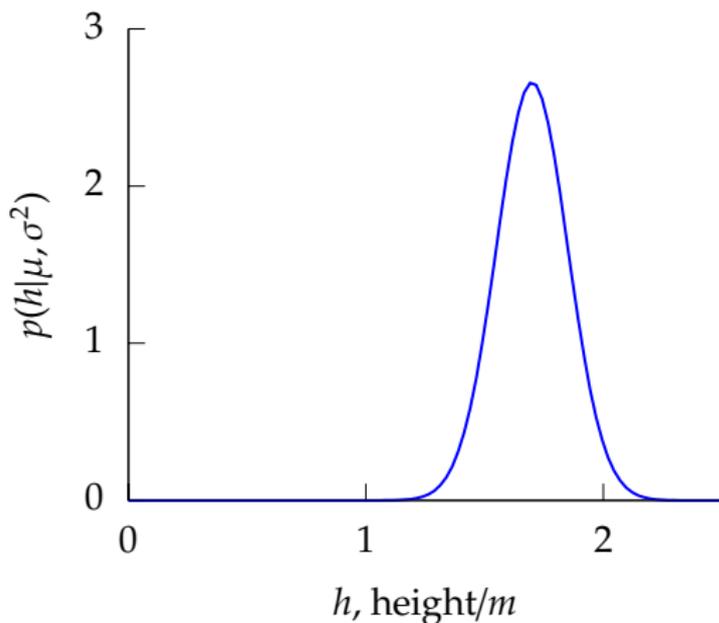
The Gaussian Density

- ▶ Perhaps the most common probability density.

$$p(y|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - \mu)^2}{2\sigma^2}\right)$$
$$\triangleq \mathcal{N}(y|\mu, \sigma^2)$$

- ▶ The Gaussian density.

Gaussian Density



The Gaussian PDF with $\mu = 1.7$ and variance $\sigma^2 = 0.0225$. Mean shown as red line. It could represent the heights of a population of students.

Gaussian Density

$$\mathcal{N}(y|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - \mu)^2}{2\sigma^2}\right)$$

σ^2 is the variance of the density and μ is the mean.

Two Important Gaussian Properties

Sum of Gaussians

- ▶ Sum of Gaussian variables is also Gaussian.

$$y_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$$

Two Important Gaussian Properties

Sum of Gaussians

- ▶ Sum of Gaussian variables is also Gaussian.

$$y_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$$

And the sum is distributed as

$$\sum_{i=1}^n y_i \sim \mathcal{N}\left(\sum_{i=1}^n \mu_i, \sum_{i=1}^n \sigma_i^2\right)$$

Two Important Gaussian Properties

Sum of Gaussians

- ▶ Sum of Gaussian variables is also Gaussian.

$$y_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$$

And the sum is distributed as

$$\sum_{i=1}^n y_i \sim \mathcal{N}\left(\sum_{i=1}^n \mu_i, \sum_{i=1}^n \sigma_i^2\right)$$

(*Aside:* As sum increases, sum of non-Gaussian, finite variance variables is also Gaussian [central limit theorem].)

Two Important Gaussian Properties

Sum of Gaussians

- ▶ Sum of Gaussian variables is also Gaussian.

$$y_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$$

And the sum is distributed as

$$\sum_{i=1}^n y_i \sim \mathcal{N}\left(\sum_{i=1}^n \mu_i, \sum_{i=1}^n \sigma_i^2\right)$$

(*Aside:* As sum increases, sum of non-Gaussian, finite variance variables is also Gaussian [central limit theorem].)

Two Important Gaussian Properties

Scaling a Gaussian

- ▶ Scaling a Gaussian leads to a Gaussian.

Two Important Gaussian Properties

Scaling a Gaussian

- ▶ Scaling a Gaussian leads to a Gaussian.

$$y \sim \mathcal{N}(\mu, \sigma^2)$$

Two Important Gaussian Properties

Scaling a Gaussian

- ▶ Scaling a Gaussian leads to a Gaussian.

$$y \sim \mathcal{N}(\mu, \sigma^2)$$

And the scaled density is distributed as

$$wy \sim \mathcal{N}(w\mu, w^2\sigma^2)$$

A Probabilistic Process

- ▶ Set the mean of Gaussian to be a function.

$$p(y_i|x_i) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - f(x_i))^2}{2\sigma^2}\right).$$

- ▶ This gives us a 'noisy function'.
- ▶ This is known as a process.

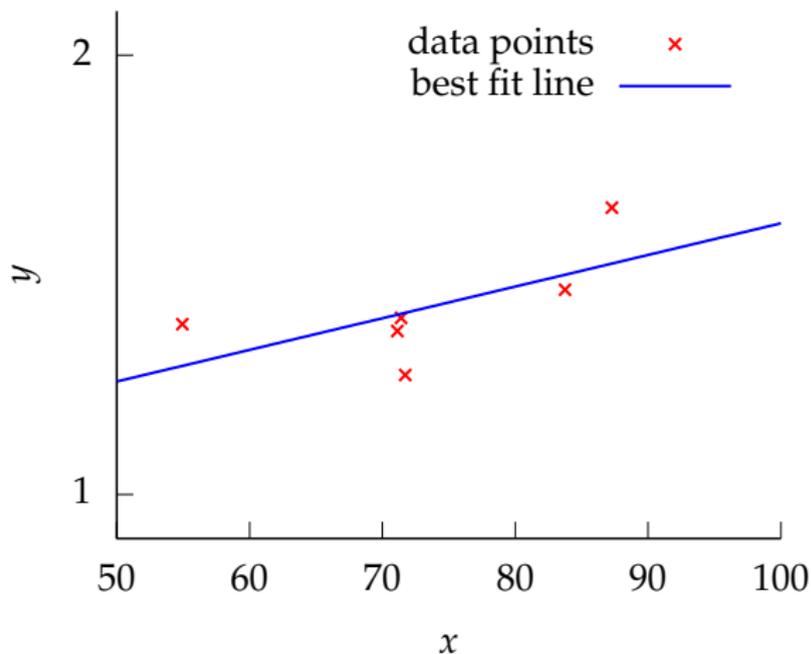
Height as a Function of Weight

- ▶ In the standard Gaussian, parametrized by mean and variance.
- ▶ Make the mean a linear function of an *input*.
- ▶ This leads to a regression model.

$$y_i = f(x_i) + \epsilon_i,$$
$$\epsilon_i \sim \mathcal{N}(0, \sigma^2).$$

- ▶ Assume y_i is height and x_i is weight.

Linear Function



A linear regression between x and y .

Data Point Likelihood

- ▶ Likelihood of an individual data point

$$p(y_i|x_i, m, c) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - mx_i - c)^2}{2\sigma^2}\right).$$

- ▶ Parameters are gradient, m , offset, c of the function and noise variance σ^2 .

Data Set Likelihood

- ▶ If the noise, ϵ_i is sampled independently for each data point.
- ▶ Each data point is independent (given m and c).
- ▶ For independent variables:

$$p(\mathbf{y}) = \prod_{i=1}^n p(y_i)$$

Data Set Likelihood

- ▶ If the noise, ϵ_i is sampled independently for each data point.
- ▶ Each data point is independent (given m and c).
- ▶ For independent variables:

$$p(\mathbf{y}|\mathbf{x}, m, c) = \prod_{i=1}^n p(y_i|x_i, m, c)$$

Data Set Likelihood

- ▶ If the noise, ϵ_i is sampled independently for each data point.
- ▶ Each data point is independent (given m and c).
- ▶ For independent variables:

$$p(\mathbf{y}|\mathbf{x}, m, c) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - mx_i - c)^2}{2\sigma^2}\right).$$

Data Set Likelihood

- ▶ If the noise, ϵ_i is sampled independently for each data point.
- ▶ Each data point is independent (given m and c).
- ▶ For independent variables:

$$p(\mathbf{y}|\mathbf{x}, m, c) = \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} \exp\left(-\frac{\sum_{i=1}^n (y_i - mx_i - c)^2}{2\sigma^2}\right).$$

Log Likelihood Function

- ▶ Normally work with the log likelihood:

$$L(m, c, \sigma^2) = -\frac{n}{2} \log 2\pi - \frac{n}{2} \log \sigma^2 - \sum_{i=1}^n \frac{(y_i - mx_i - c)^2}{2\sigma^2}.$$

Consistency of Maximum Likelihood

- ▶ If data was really generated according to probability we specified.
- ▶ Correct parameters will be recovered in limit as $n \rightarrow \infty$.
- ▶ This can be proven through sample based approximations (law of large numbers) of “KL divergences”.
- ▶ Mainstay of classical statistics.

Probabilistic Interpretation of the Error Function

- ▶ Probabilistic Interpretation for Error Function is Negative Log Likelihood.
- ▶ *Minimizing* error function is equivalent to *maximizing* log likelihood.
- ▶ Maximizing *log likelihood* is equivalent to maximizing the *likelihood* because log is monotonic.
- ▶ Probabilistic interpretation: Minimizing error function is equivalent to maximum likelihood with respect to parameters.

Error Function

- ▶ Negative log likelihood is the error function leading to an error function

$$E(m, c, \sigma^2) = \frac{n}{2} \log \sigma^2 + \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - mx_i - c)^2.$$

- ▶ Learning proceeds by minimizing this error function for the data set provided.

Connection: Sum of Squares Error

- ▶ Ignoring terms which don't depend on m and c gives

$$E(m, c) \propto \sum_{i=1}^n (y_i - f(x_i))^2$$

where $f(x_i) = mx_i + c$.

- ▶ This is known as the *sum of squares* error function.
- ▶ Commonly used and is closely associated with the Gaussian likelihood.

Mathematical Interpretation

- ▶ What is the mathematical interpretation?
 - ▶ There is a cost function.
 - ▶ It expresses mismatch between your prediction and reality.

$$E(m, c) = \sum_{i=1}^n (y_i - mx_i - c)^2$$

- ▶ This is known as the sum of squares error.

Learning is Optimization

- ▶ Learning is minimization of the cost function.
- ▶ At the minima the gradient is zero.
- ▶ Coordinate ascent, find gradient in each coordinate and set to zero.

$$\frac{dE(m)}{dm} = -2 \sum_{i=1}^n x_i (y_i - mx_i - c)$$

Learning is Optimization

- ▶ Learning is minimization of the cost function.
- ▶ At the minima the gradient is zero.
- ▶ Coordinate ascent, find gradient in each coordinate and set to zero.

$$0 = -2 \sum_{i=1}^n x_i (y_i - mx_i - c)$$

Learning is Optimization

- ▶ Learning is minimization of the cost function.
- ▶ At the minima the gradient is zero.
- ▶ Coordinate ascent, find gradient in each coordinate and set to zero.

$$0 = -2 \sum_{i=1}^n x_i y_i + 2 \sum_{i=1}^n m x_i^2 + 2 \sum_{i=1}^n c x_i$$

Learning is Optimization

- ▶ Learning is minimization of the cost function.
- ▶ At the minima the gradient is zero.
- ▶ Coordinate ascent, find gradient in each coordinate and set to zero.

$$m = \frac{\sum_{i=1}^n (y_i - c) x_i}{\sum_{i=1}^n x_i^2}$$

Learning is Optimization

- ▶ Learning is minimization of the cost function.
- ▶ At the minima the gradient is zero.
- ▶ Coordinate ascent, find gradient in each coordinate and set to zero.

$$\frac{dE(c)}{dc} = -2 \sum_{i=1}^n (y_i - mx_i - c)$$

Learning is Optimization

- ▶ Learning is minimization of the cost function.
- ▶ At the minima the gradient is zero.
- ▶ Coordinate ascent, find gradient in each coordinate and set to zero.

$$0 = -2 \sum_{i=1}^n (y_i - mx_i - c)$$

Learning is Optimization

- ▶ Learning is minimization of the cost function.
- ▶ At the minima the gradient is zero.
- ▶ Coordinate ascent, find gradient in each coordinate and set to zero.

$$0 = -2 \sum_{i=1}^n y_i + 2 \sum_{i=1}^n mx_i + 2nc$$

Learning is Optimization

- ▶ Learning is minimization of the cost function.
- ▶ At the minima the gradient is zero.
- ▶ Coordinate ascent, find gradient in each coordinate and set to zero.

$$c = \frac{\sum_{i=1}^n (y_i - mx_i)}{n}$$

Fixed Point Updates

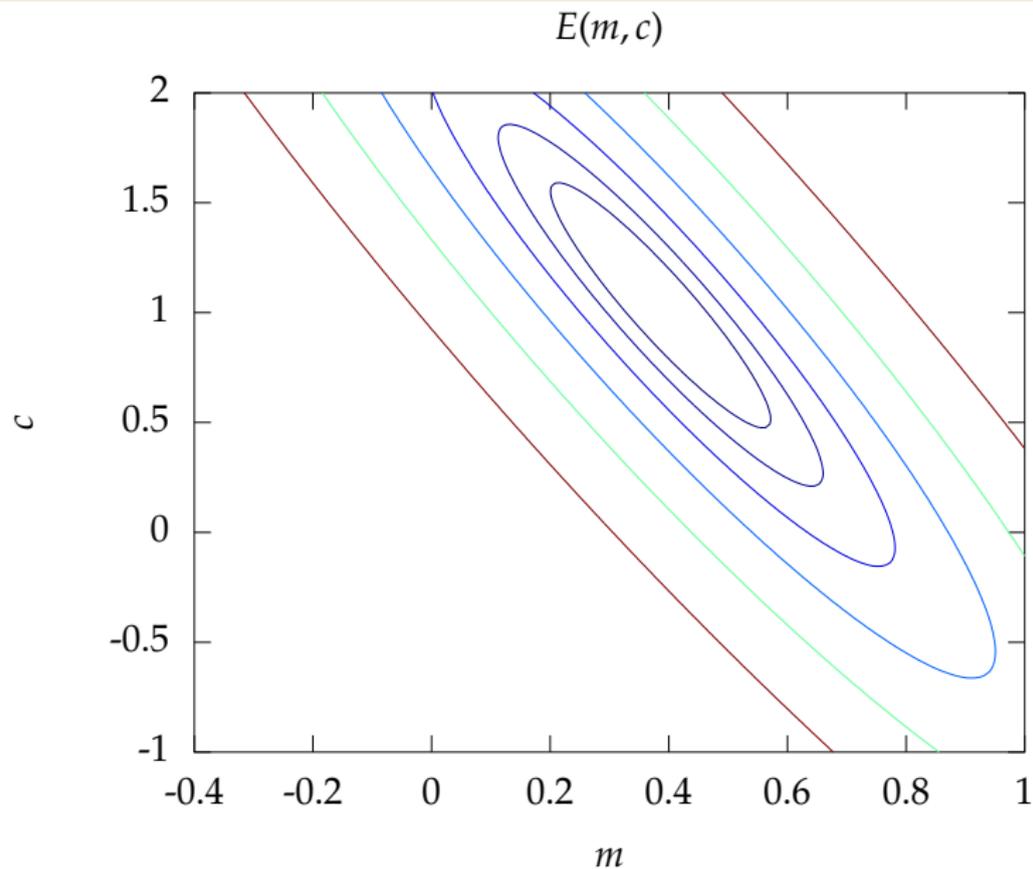
Worked example.

$$c^* = \frac{\sum_{i=1}^n (y_i - m^* x_i)}{n},$$

$$m^* = \frac{\sum_{i=1}^n x_i (y_i - c^*)}{\sum_{i=1}^n x_i^2},$$

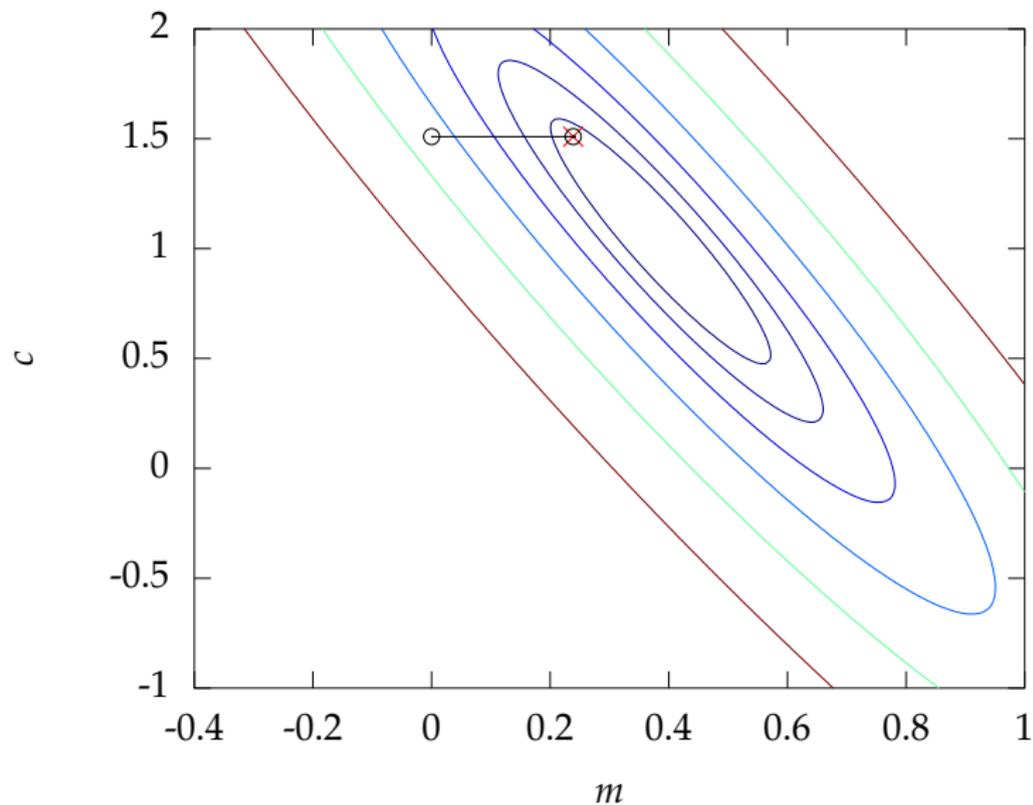
$$\sigma^{2*} = \frac{\sum_{i=1}^n (y_i - m^* x_i - c^*)^2}{n}$$

Coordinate Descent



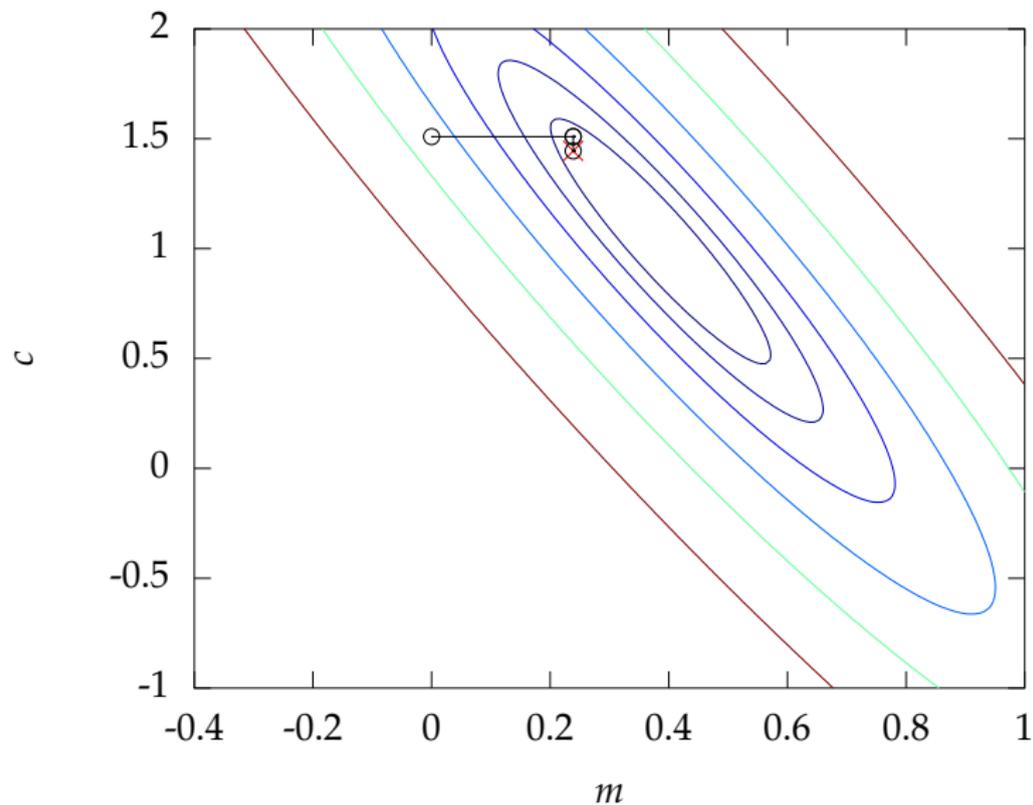
Coordinate Descent

Iteration 1



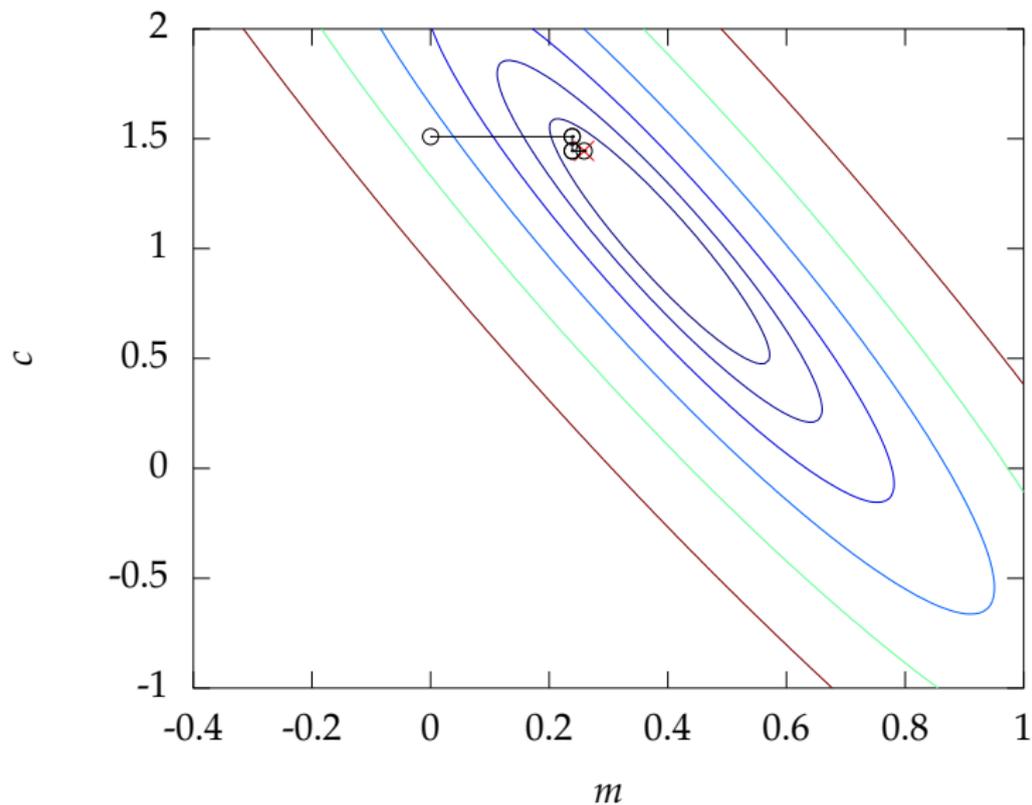
Coordinate Descent

Iteration 1



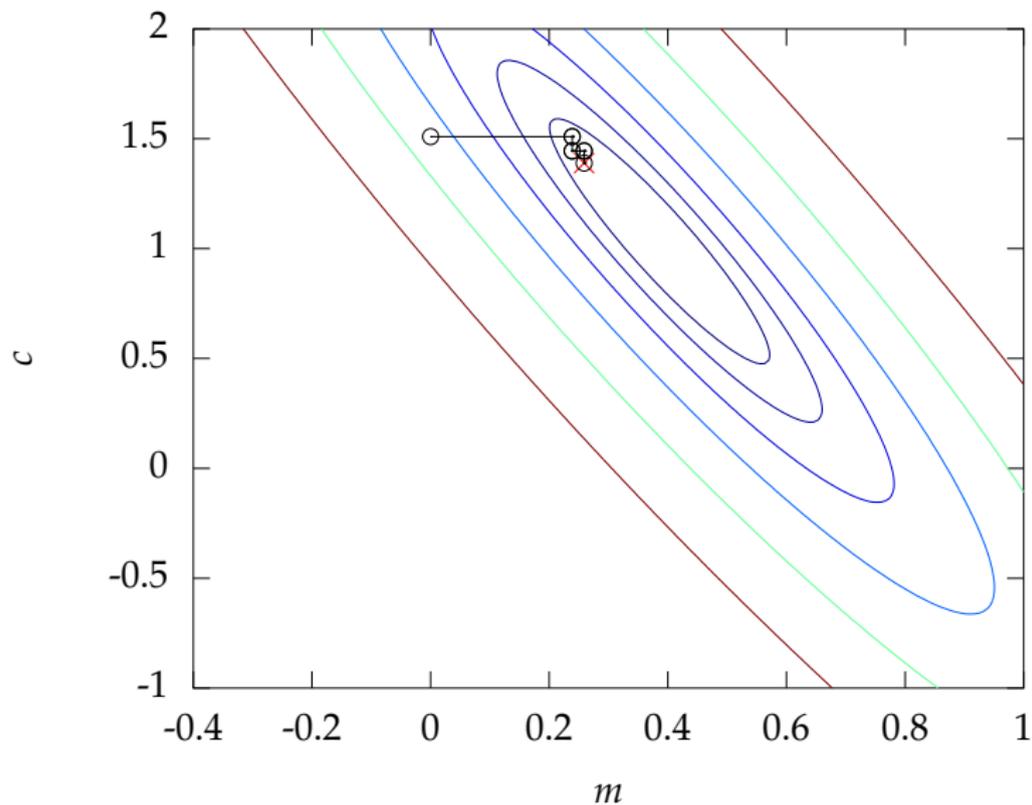
Coordinate Descent

Iteration 2



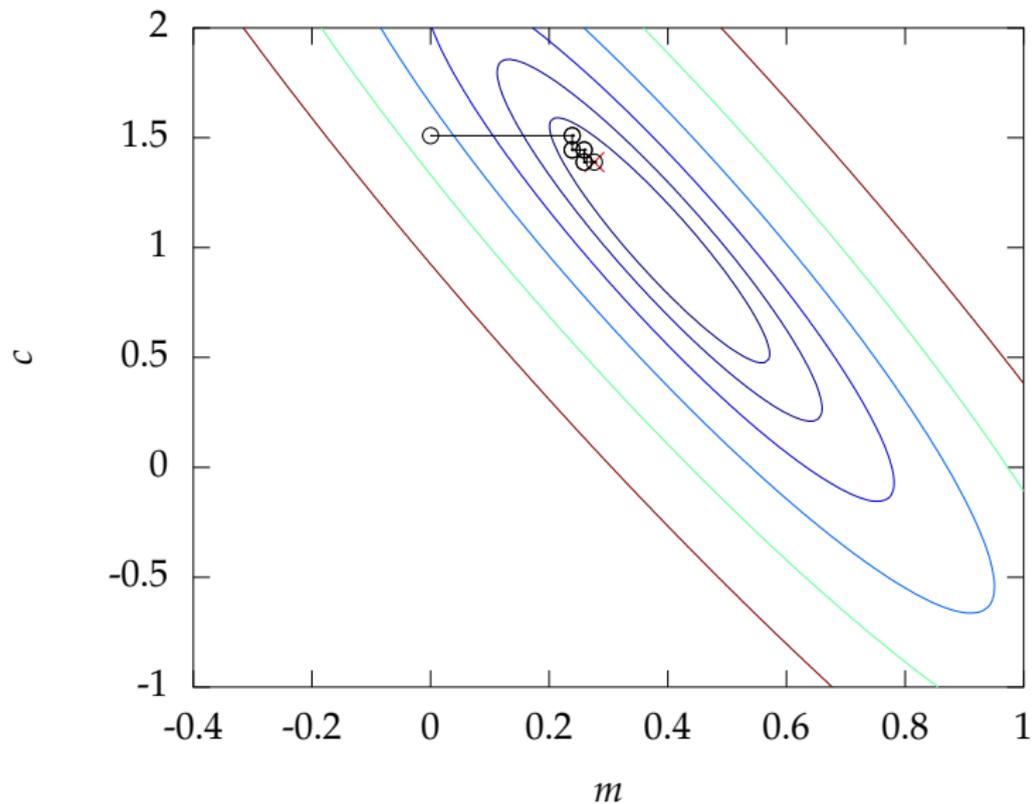
Coordinate Descent

Iteration 2



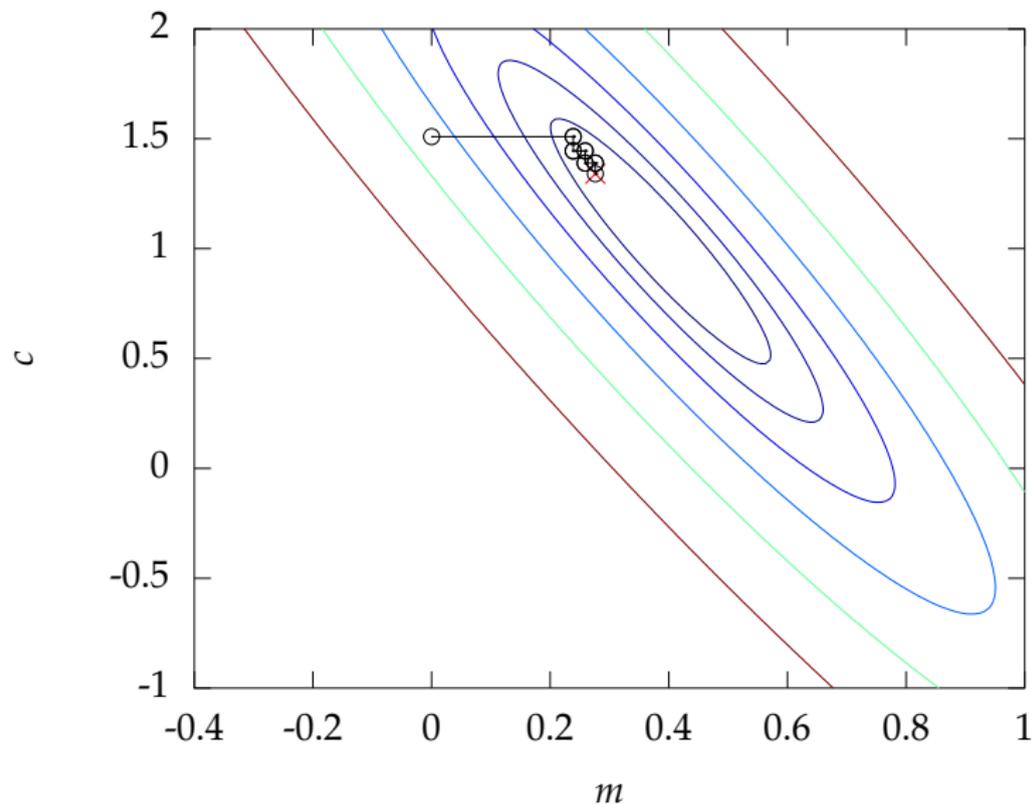
Coordinate Descent

Iteration 3



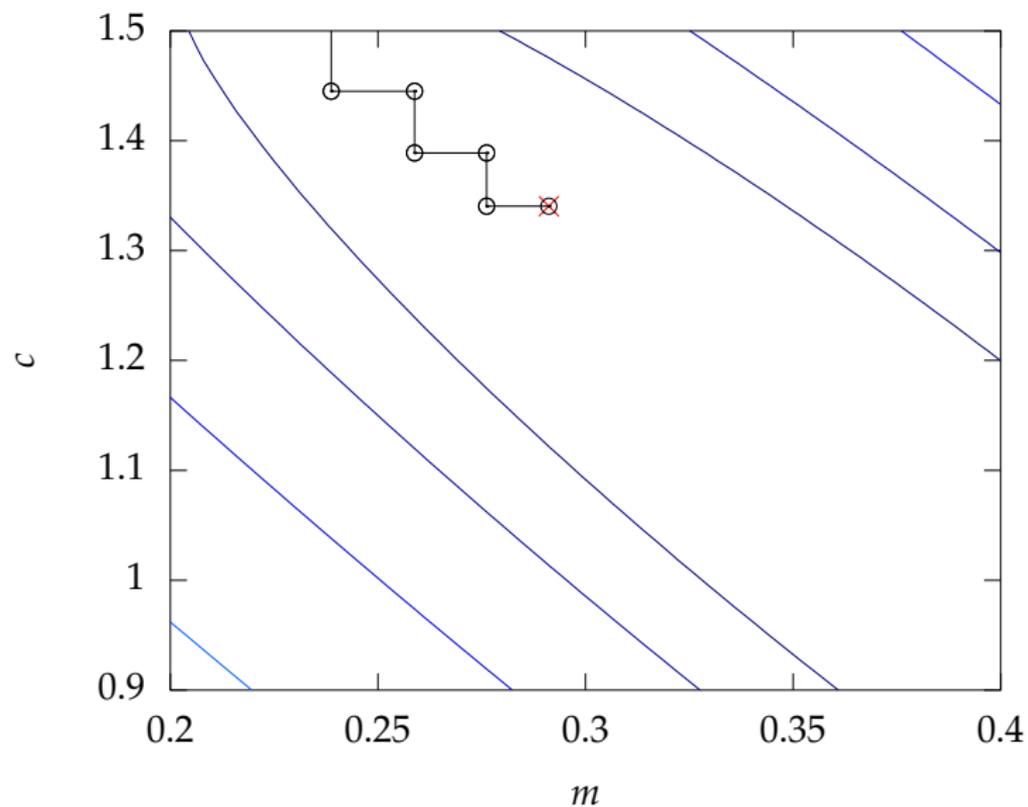
Coordinate Descent

Iteration 3



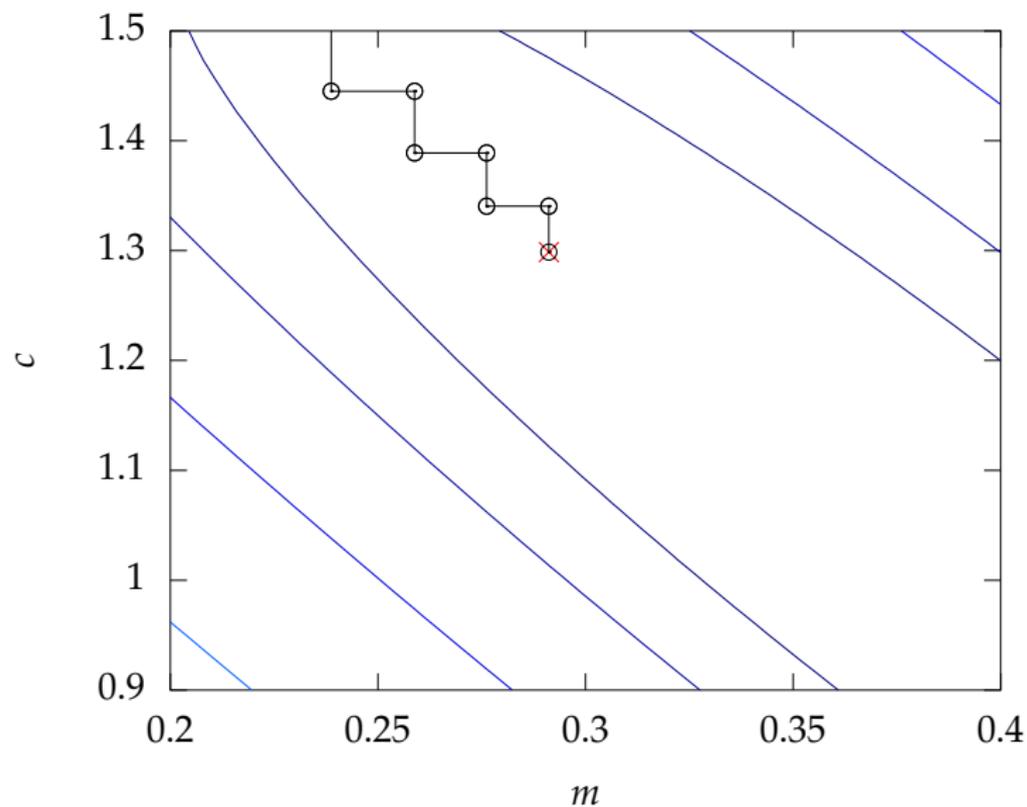
Coordinate Descent

Iteration 4



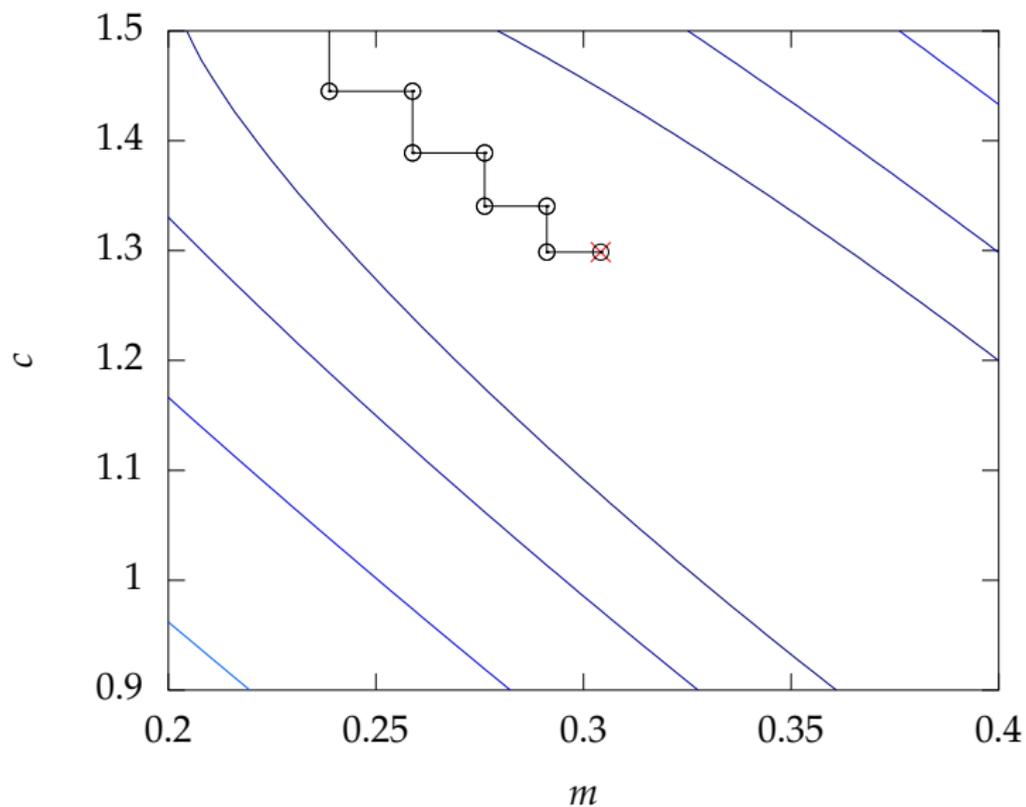
Coordinate Descent

Iteration 4



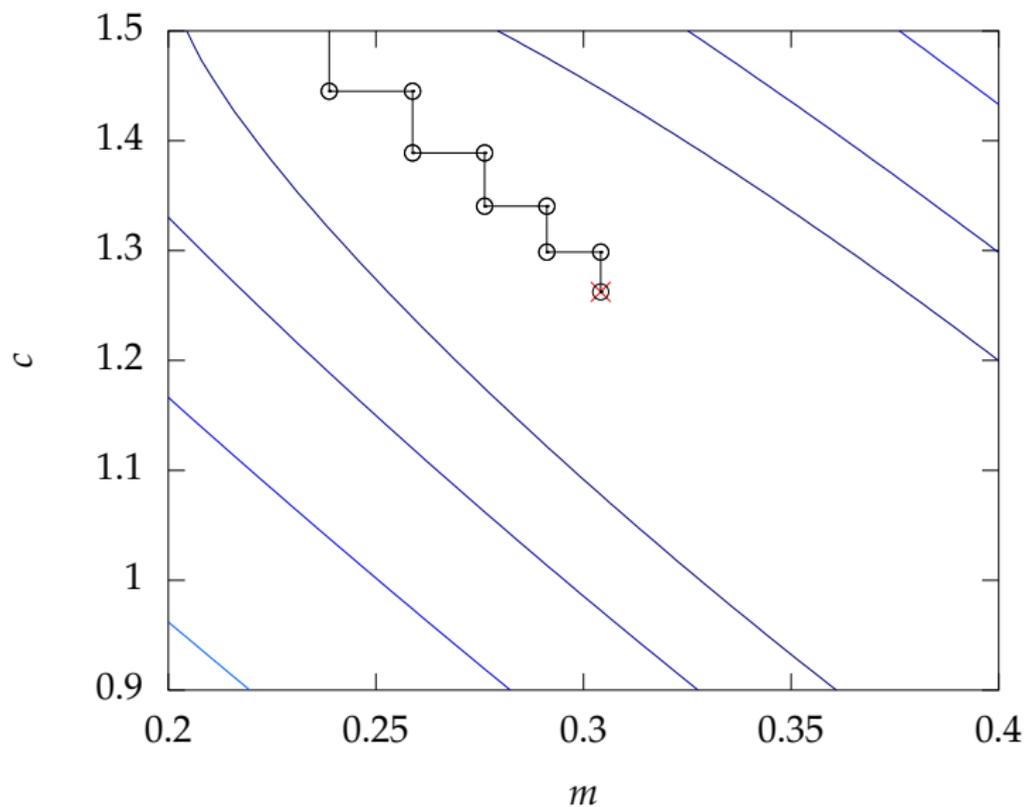
Coordinate Descent

Iteration 5



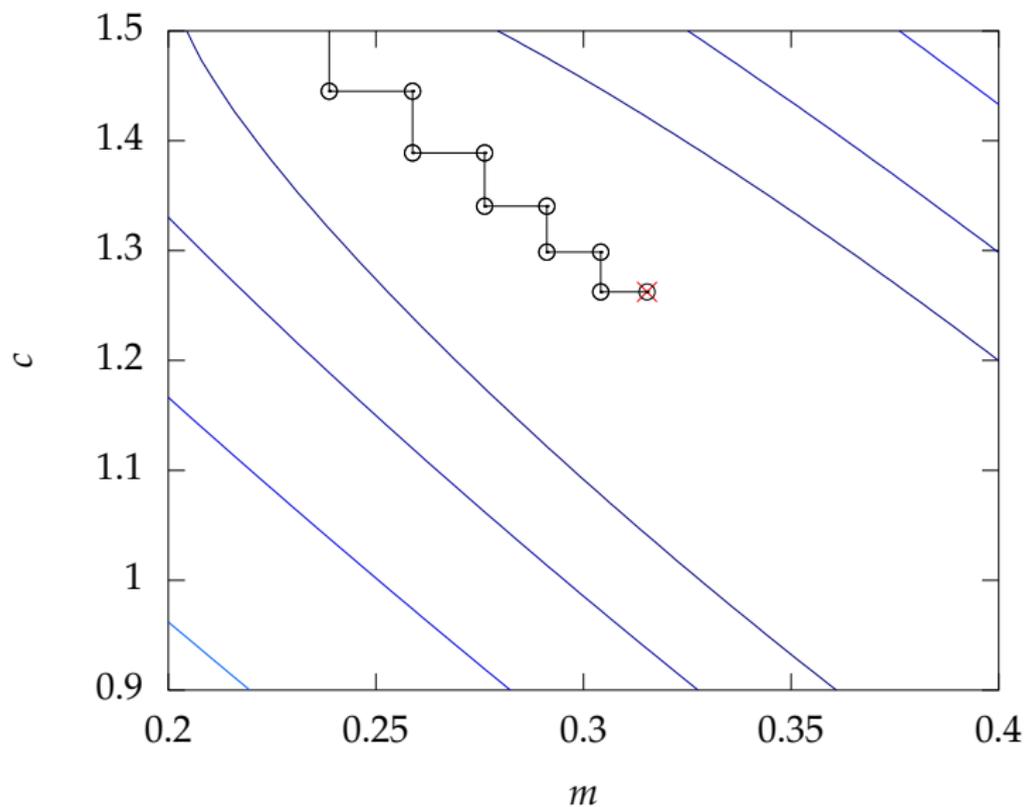
Coordinate Descent

Iteration 5



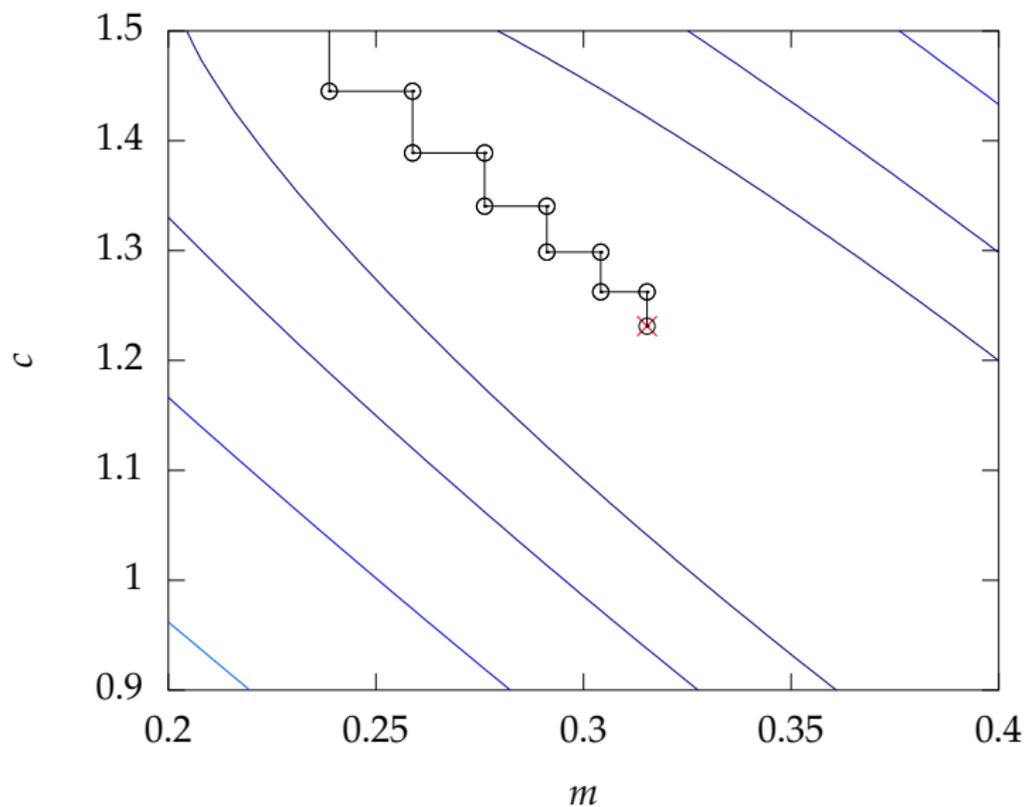
Coordinate Descent

Iteration 6



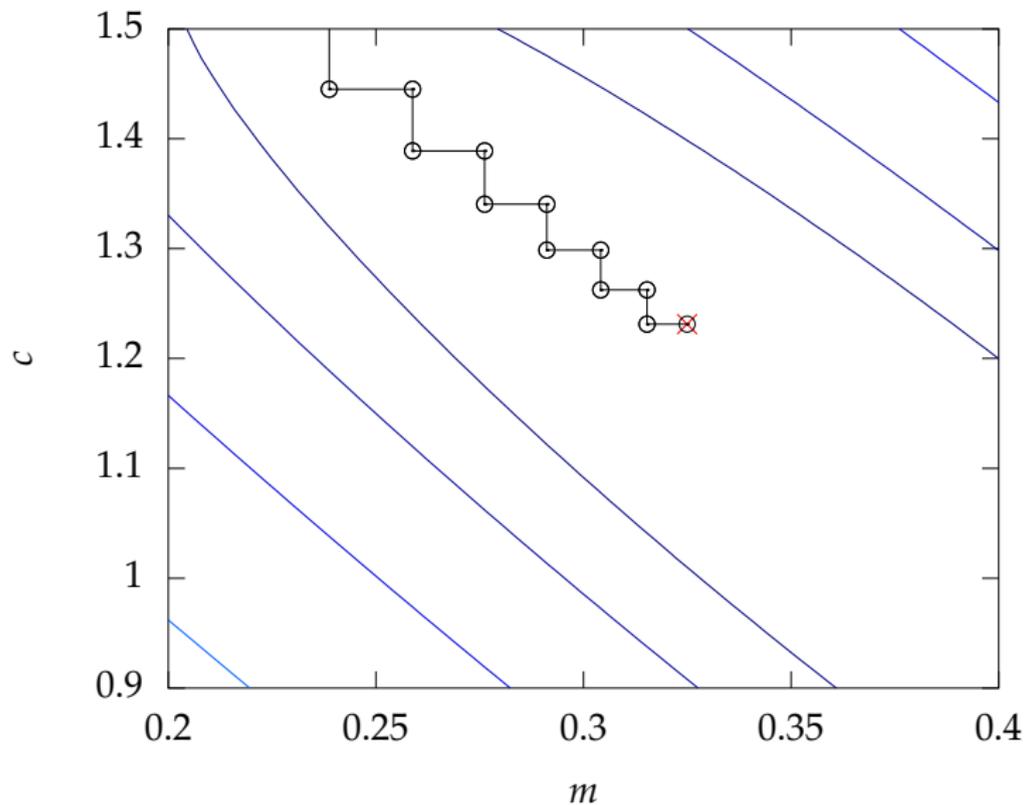
Coordinate Descent

Iteration 6



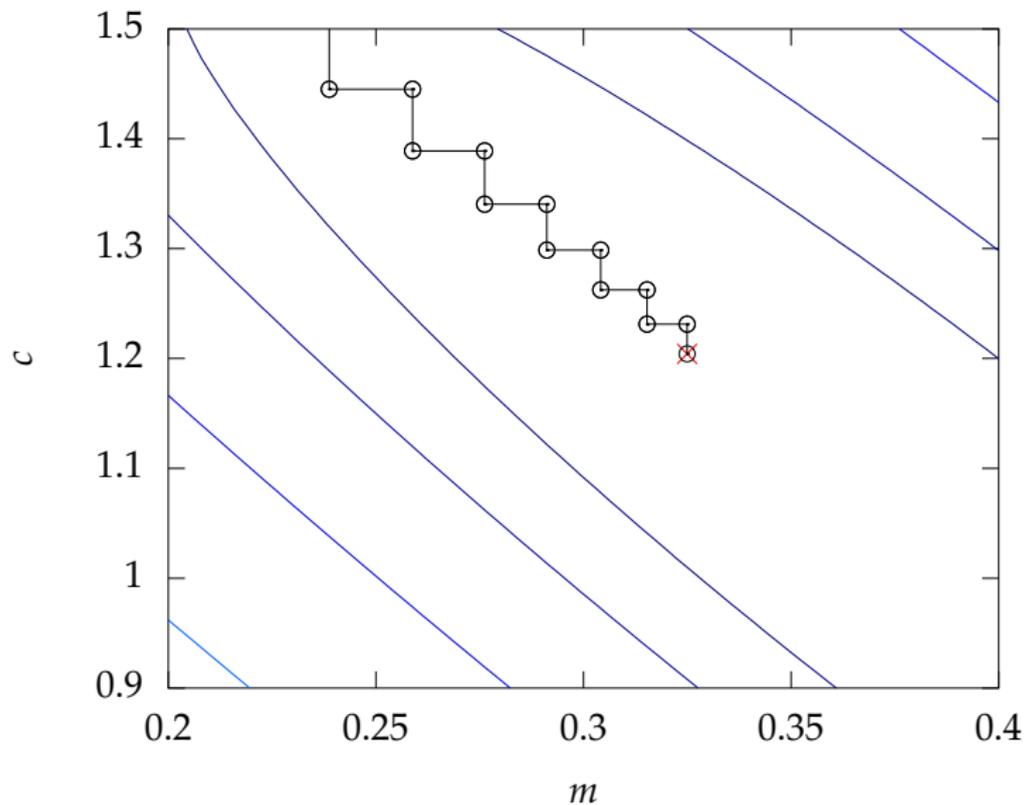
Coordinate Descent

Iteration 7



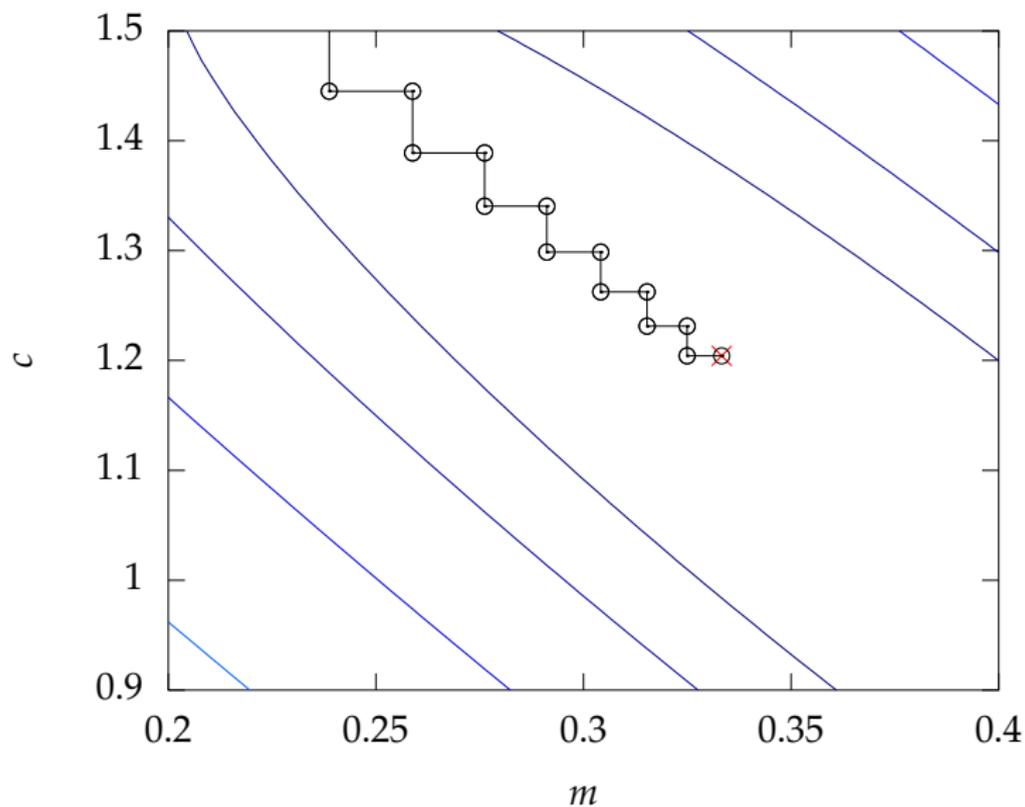
Coordinate Descent

Iteration 7



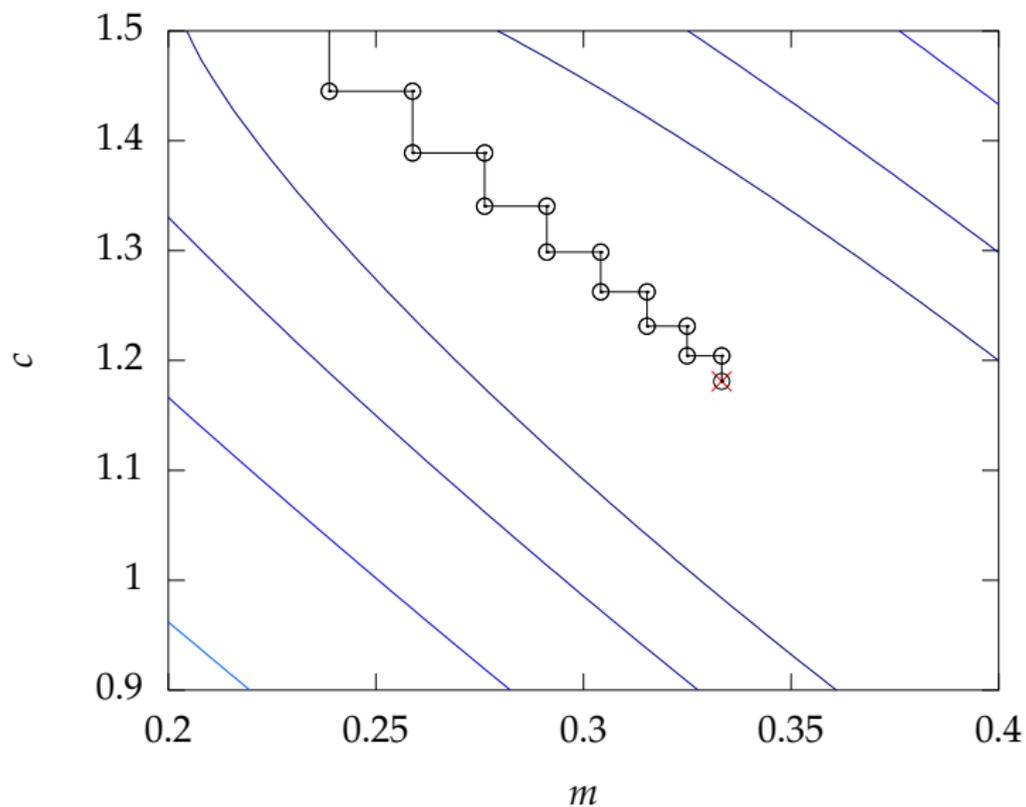
Coordinate Descent

Iteration 8



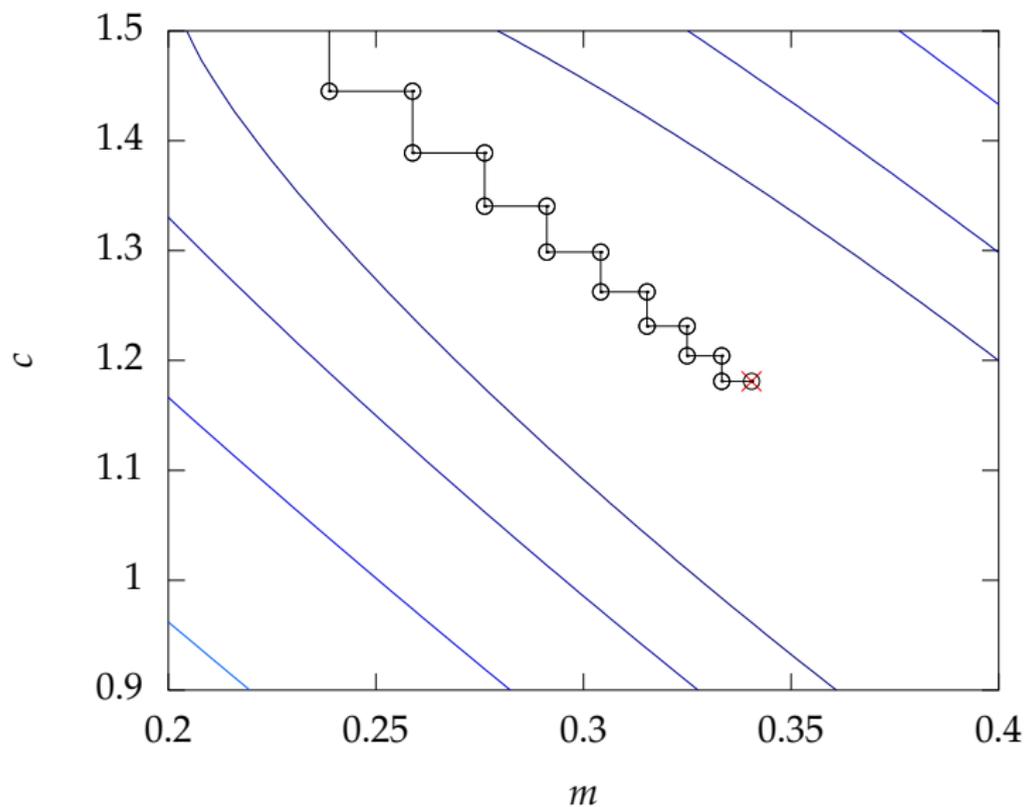
Coordinate Descent

Iteration 8



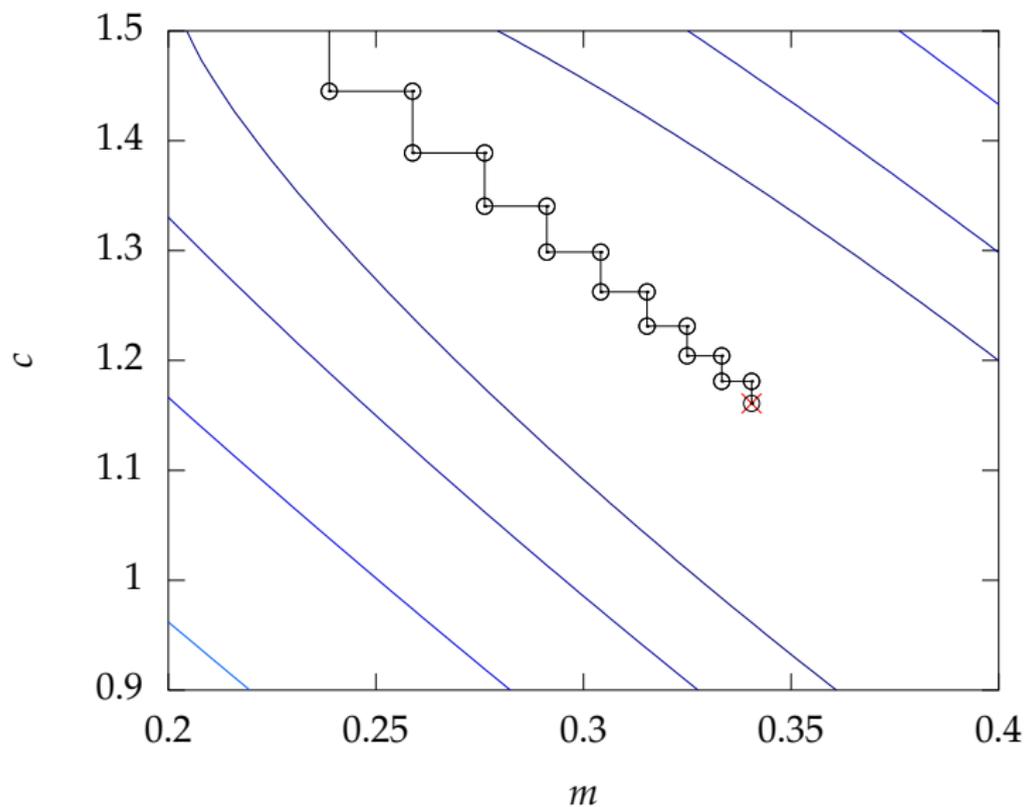
Coordinate Descent

Iteration 9



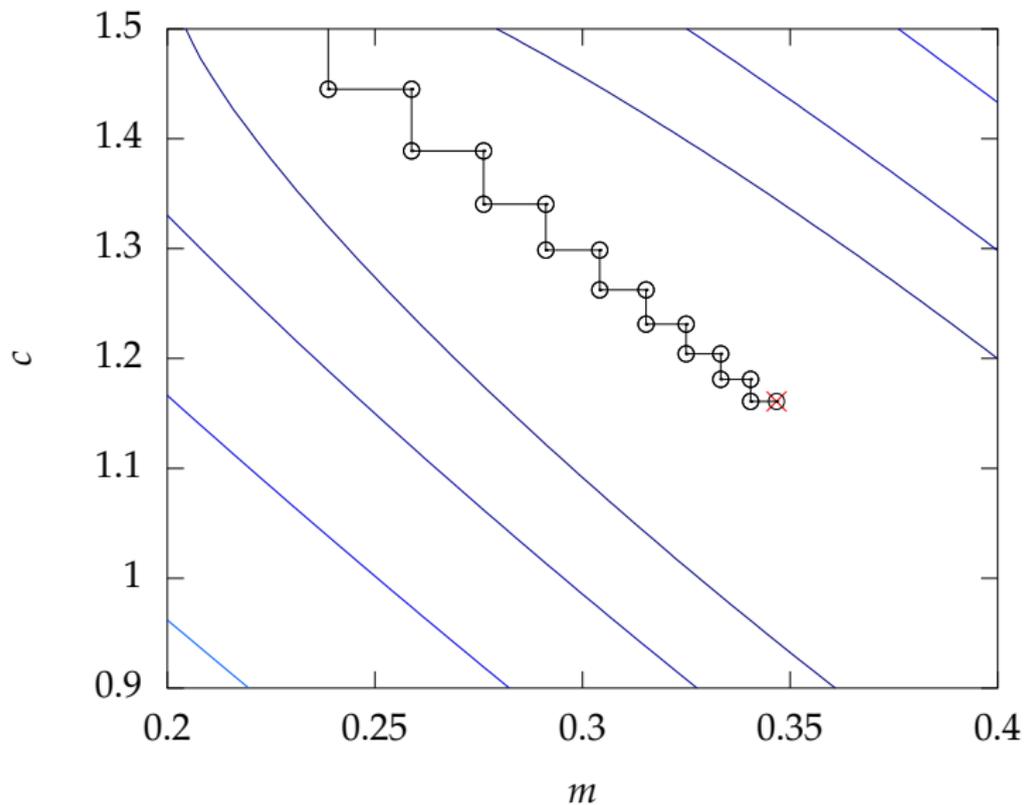
Coordinate Descent

Iteration 9



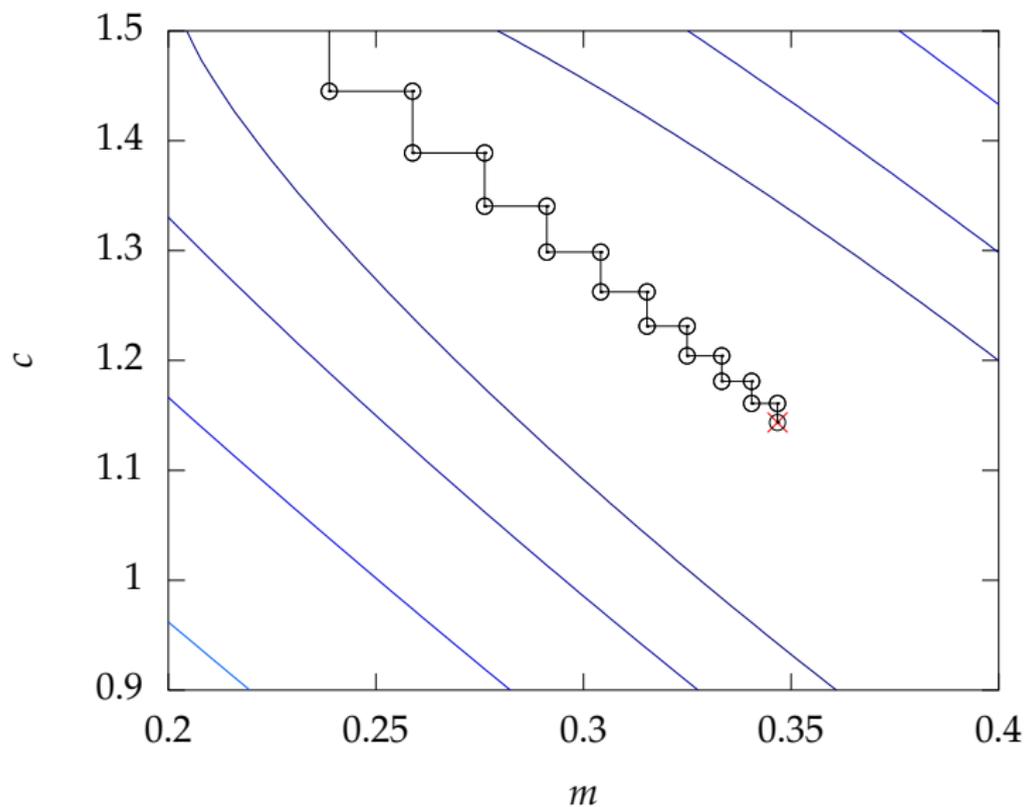
Coordinate Descent

Iteration 10



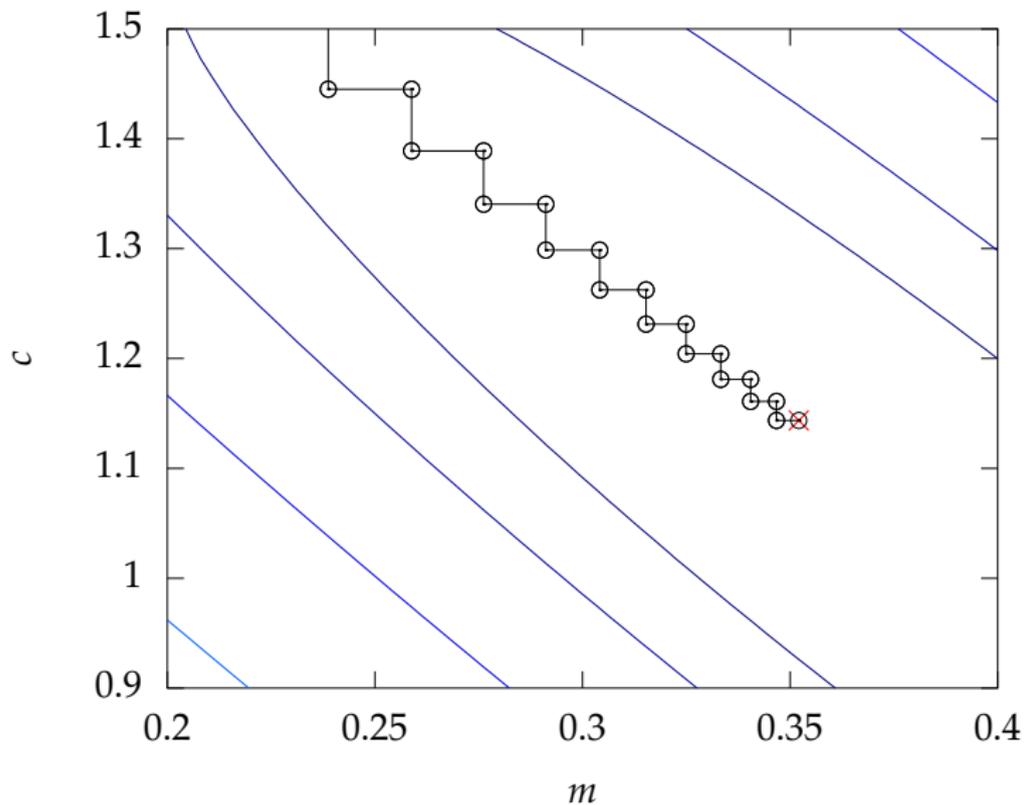
Coordinate Descent

Iteration 10



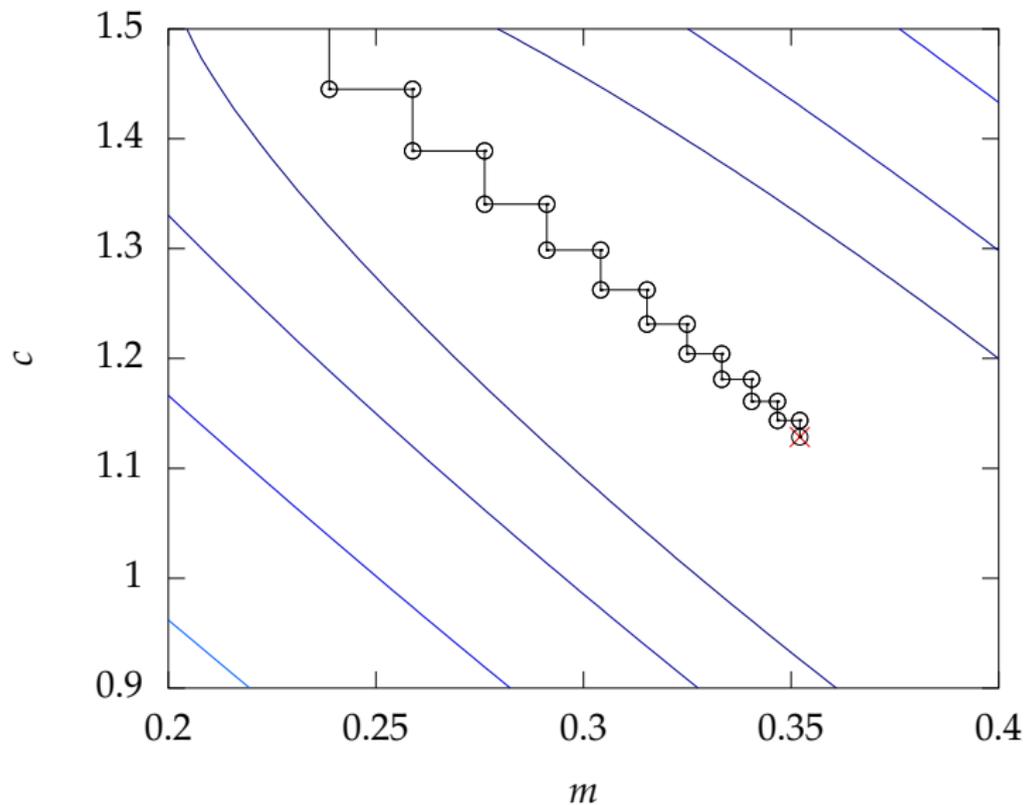
Coordinate Descent

Iteration 10



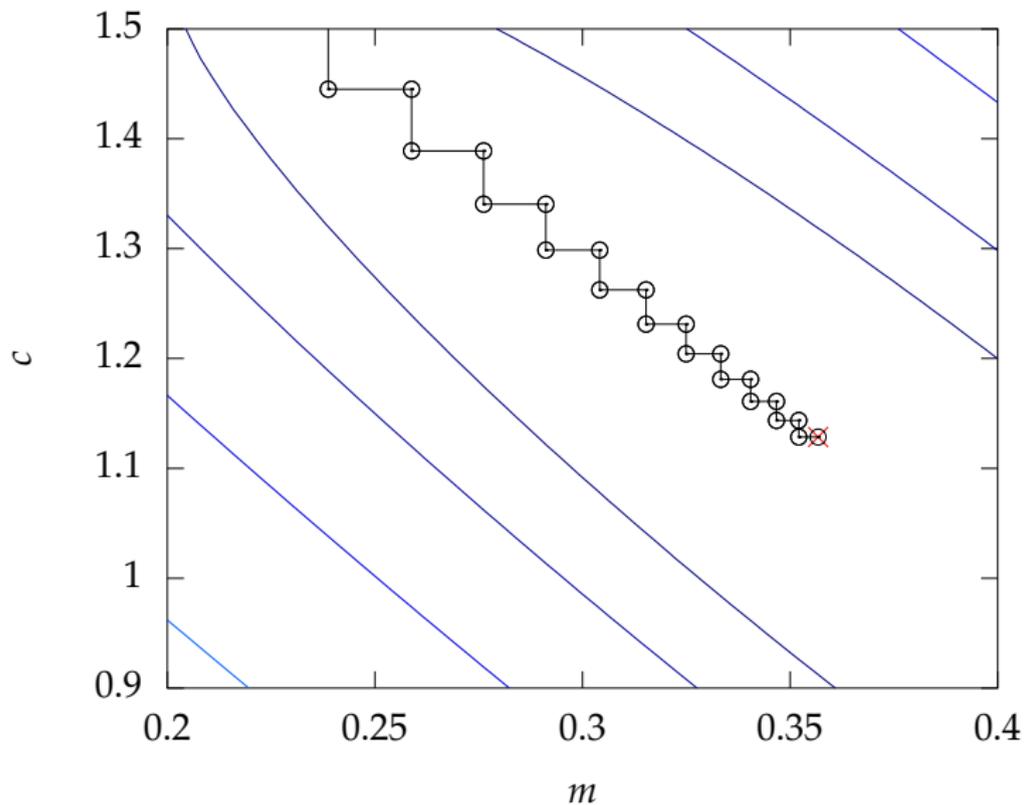
Coordinate Descent

Iteration 10



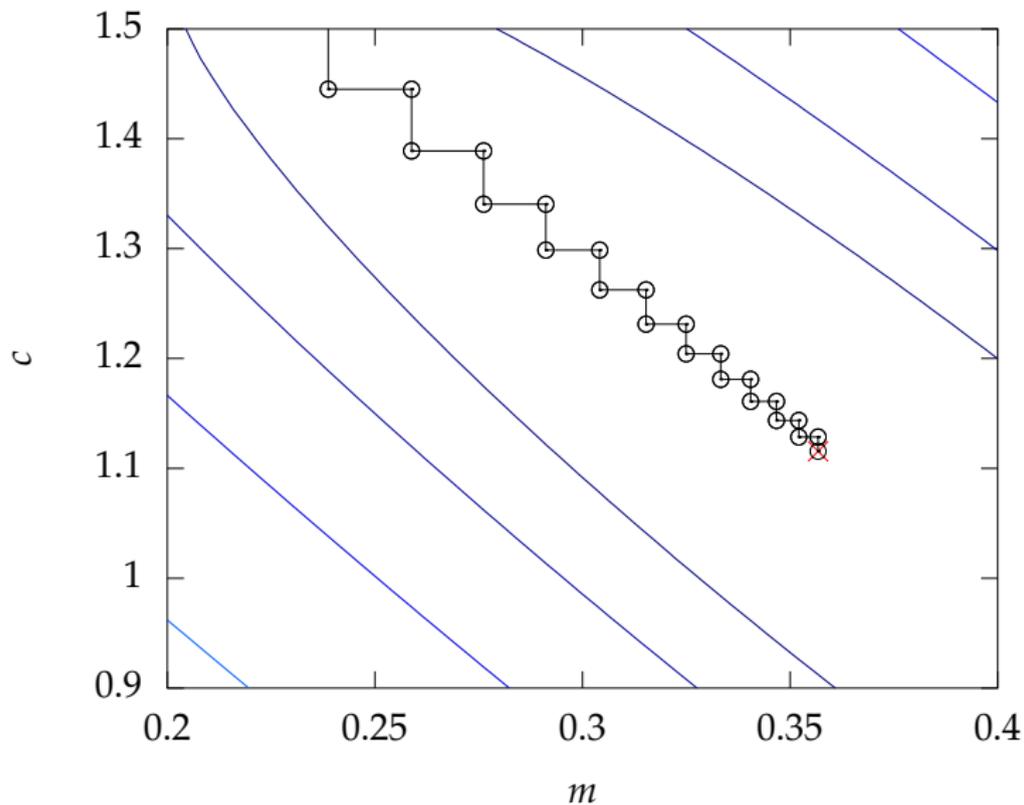
Coordinate Descent

Iteration 10



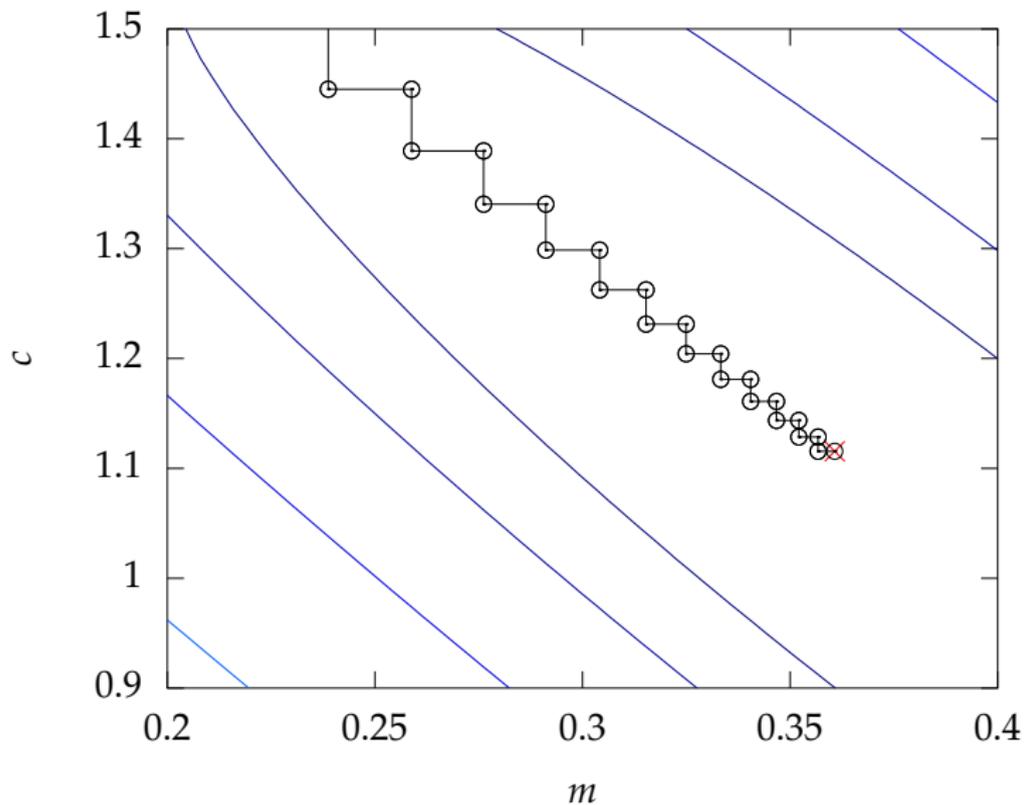
Coordinate Descent

Iteration 10



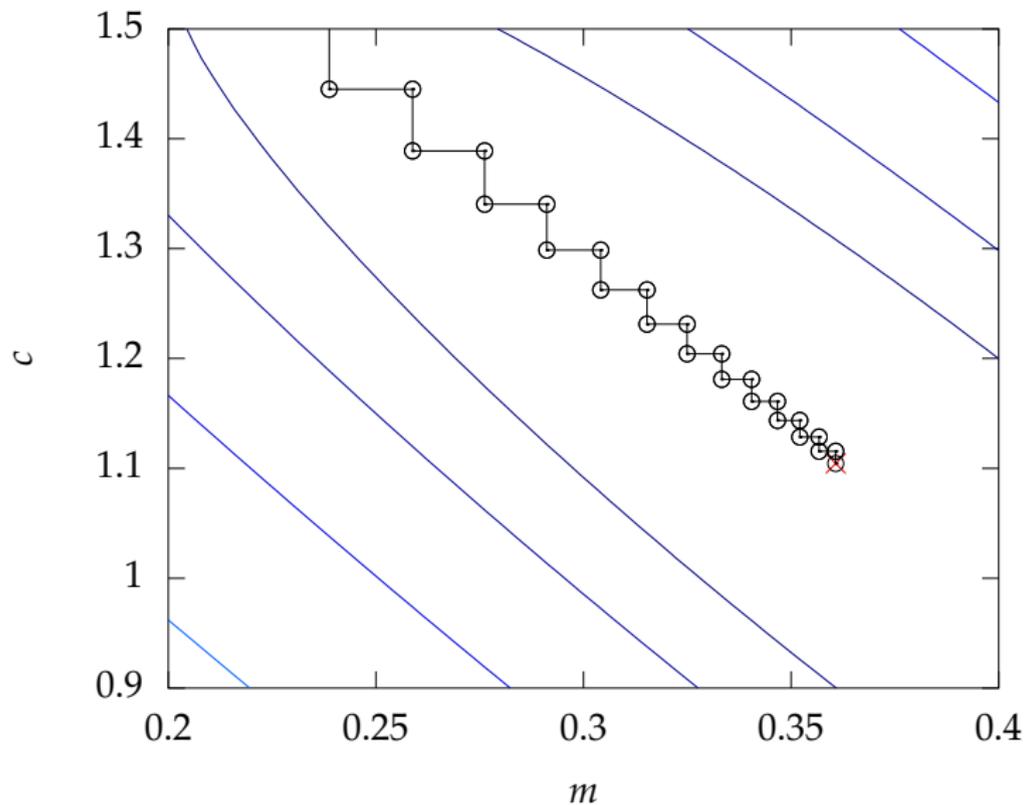
Coordinate Descent

Iteration 10



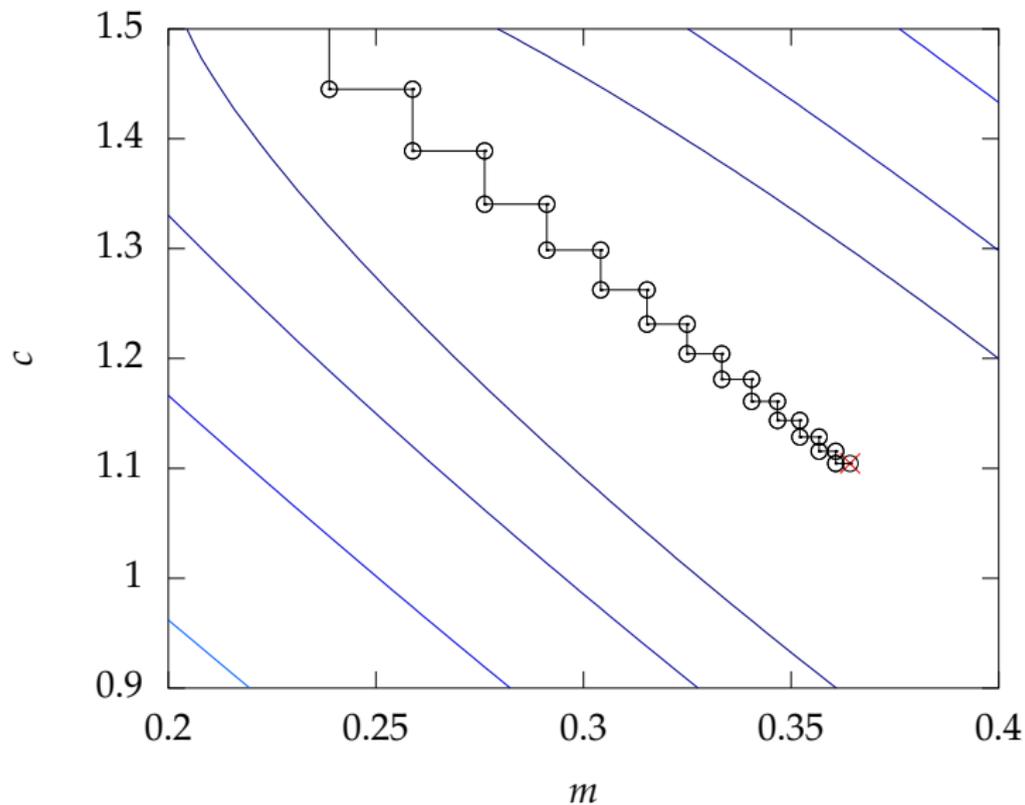
Coordinate Descent

Iteration 10



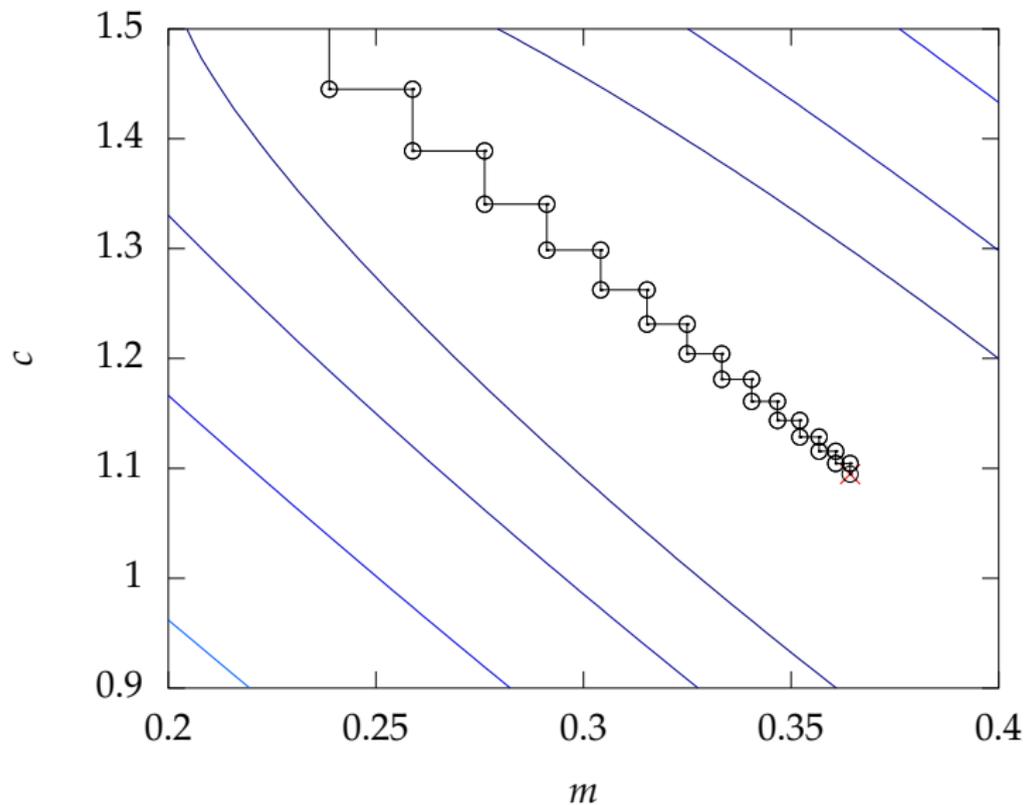
Coordinate Descent

Iteration 10



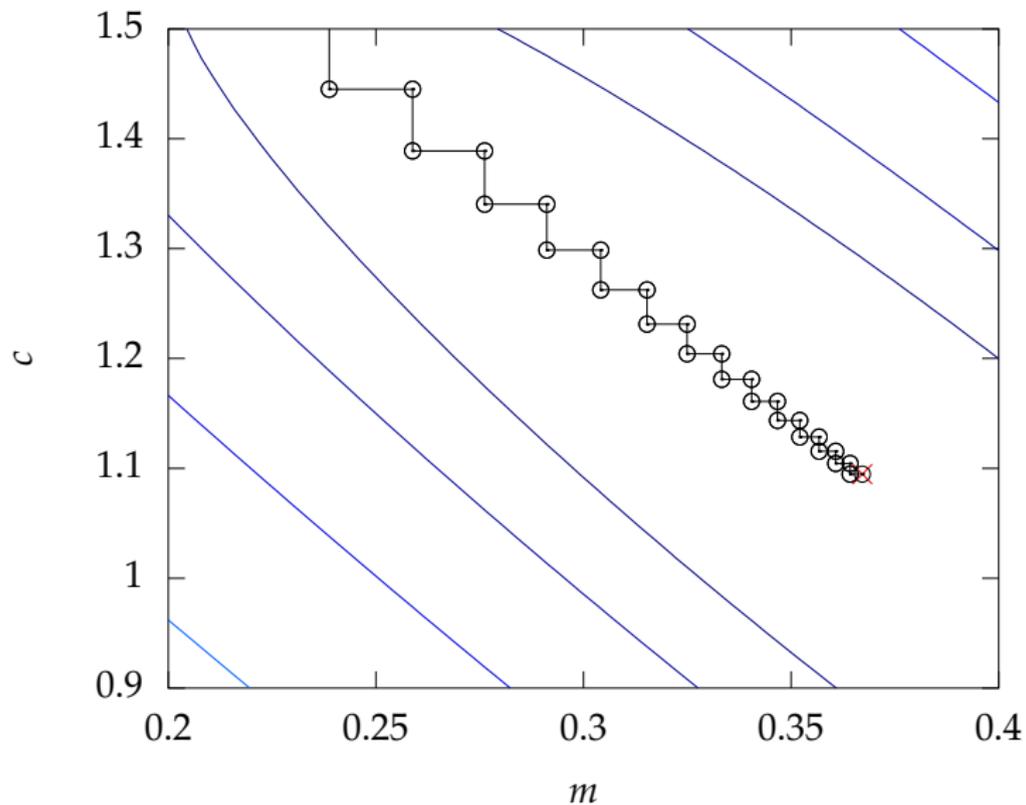
Coordinate Descent

Iteration 10



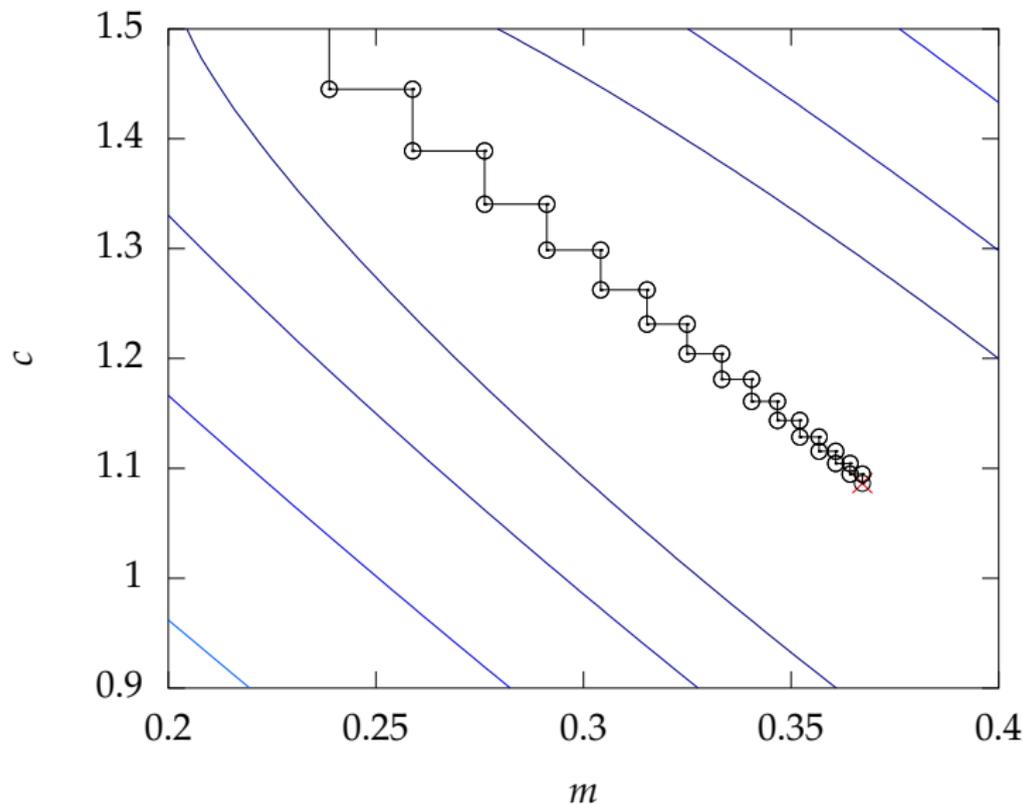
Coordinate Descent

Iteration 10



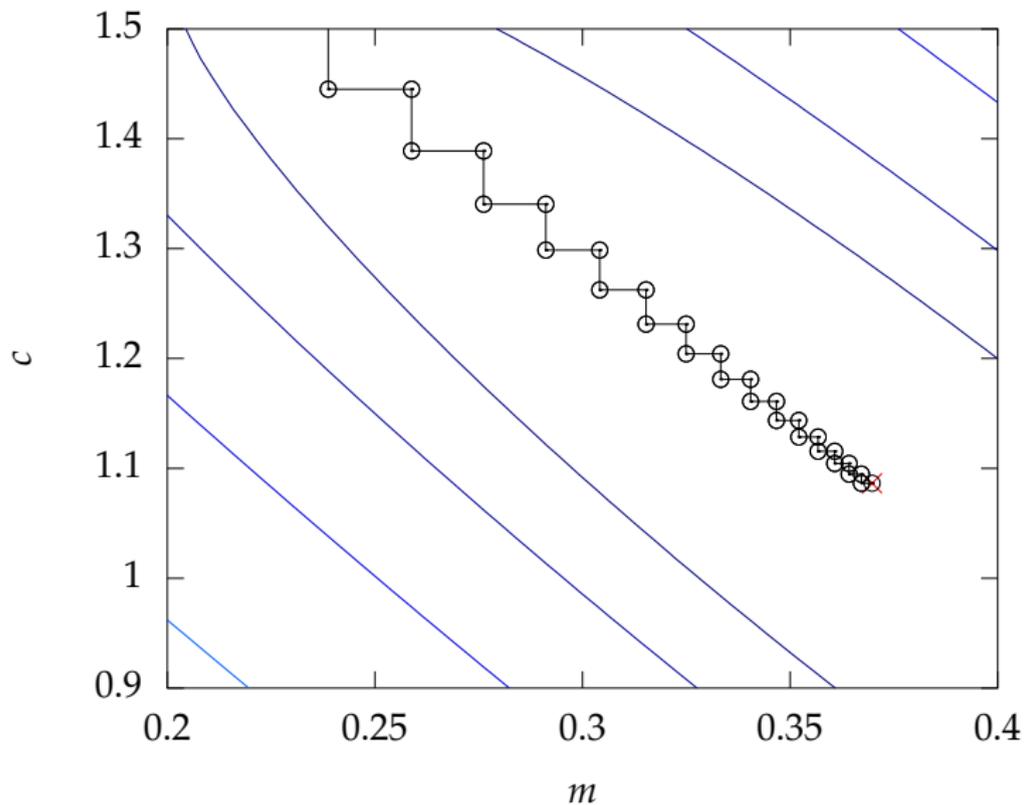
Coordinate Descent

Iteration 10



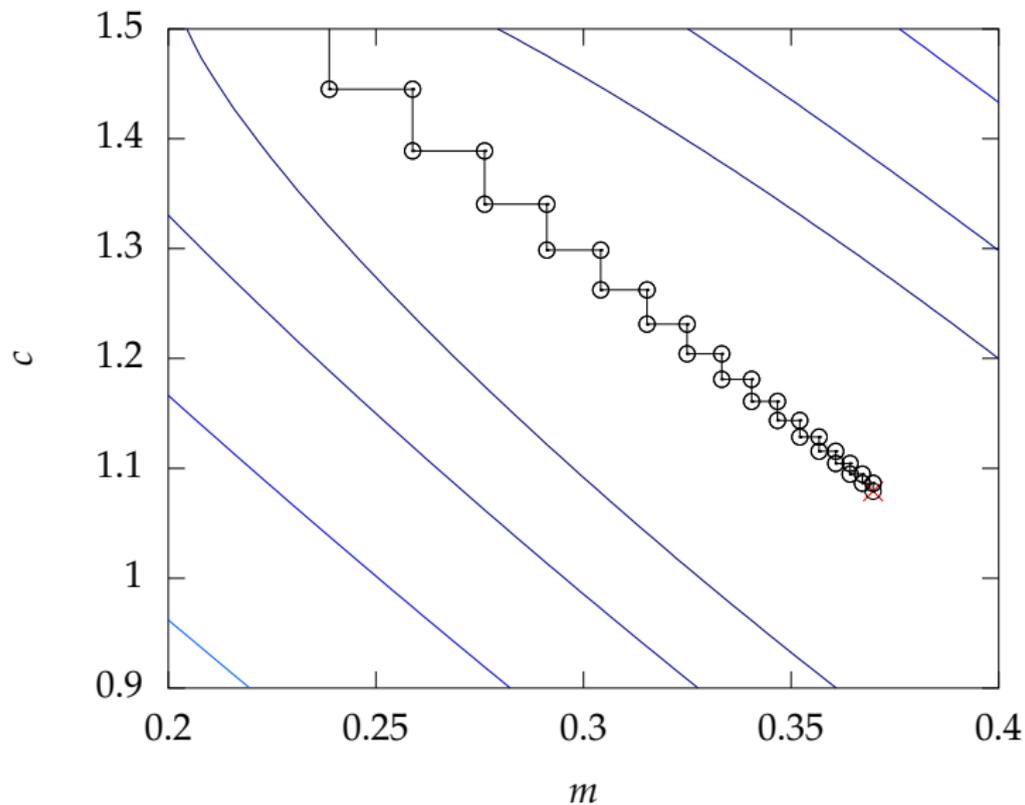
Coordinate Descent

Iteration 10



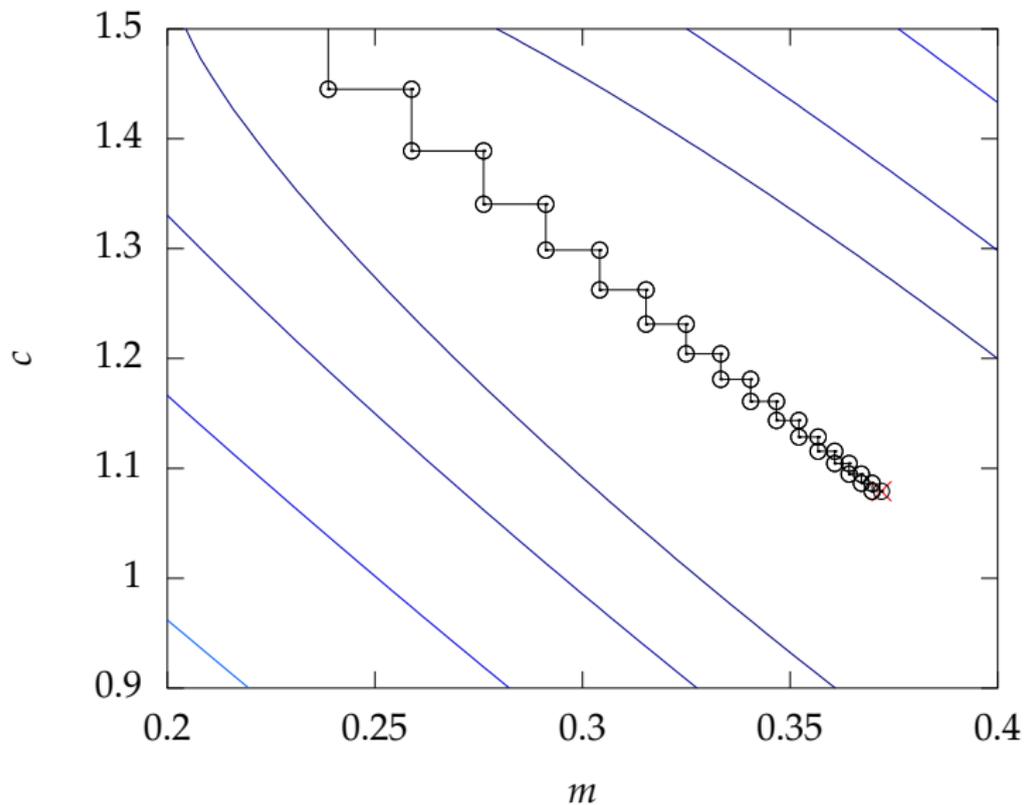
Coordinate Descent

Iteration 10



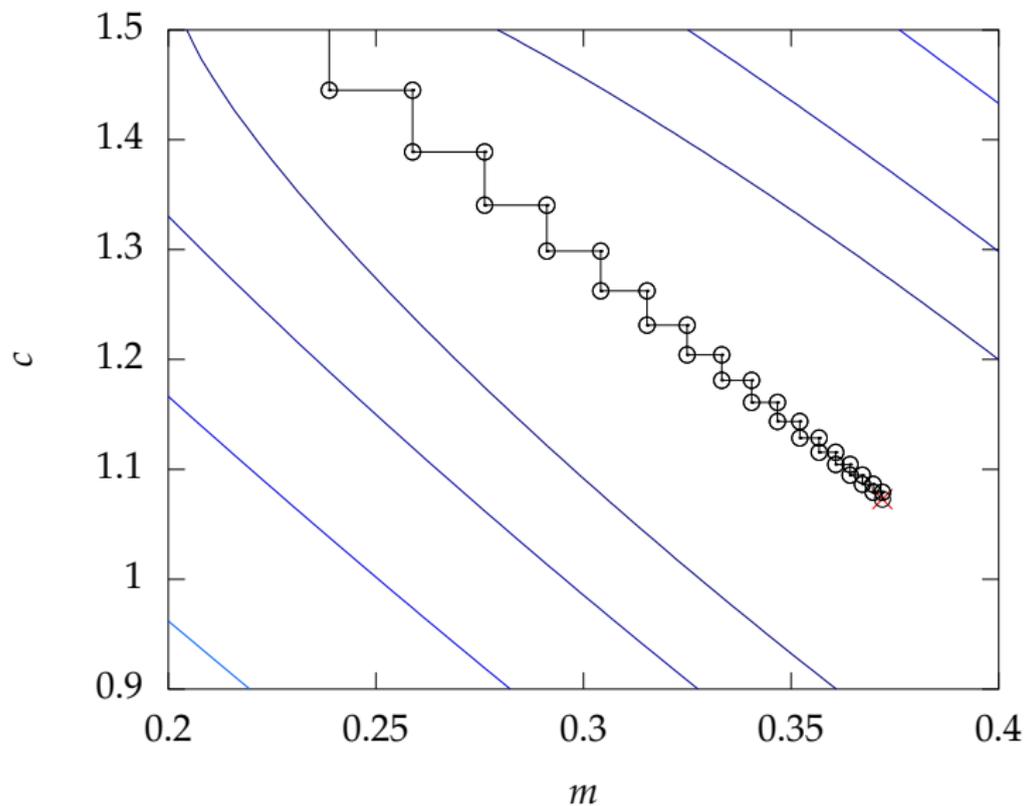
Coordinate Descent

Iteration 10



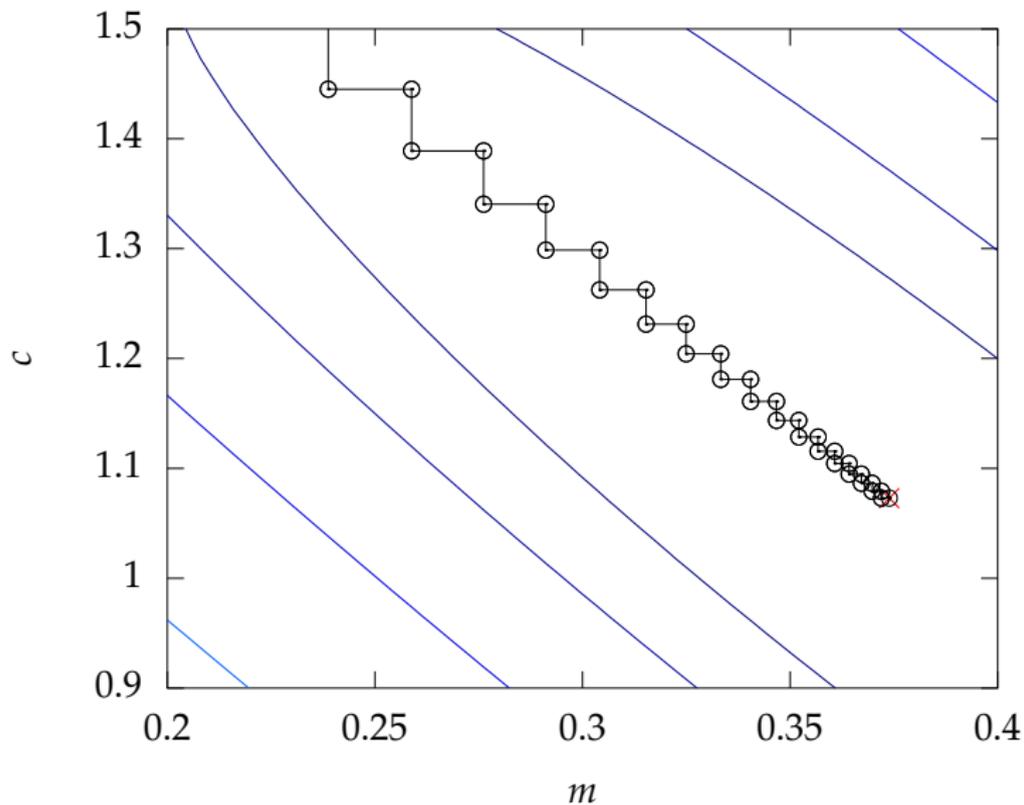
Coordinate Descent

Iteration 10



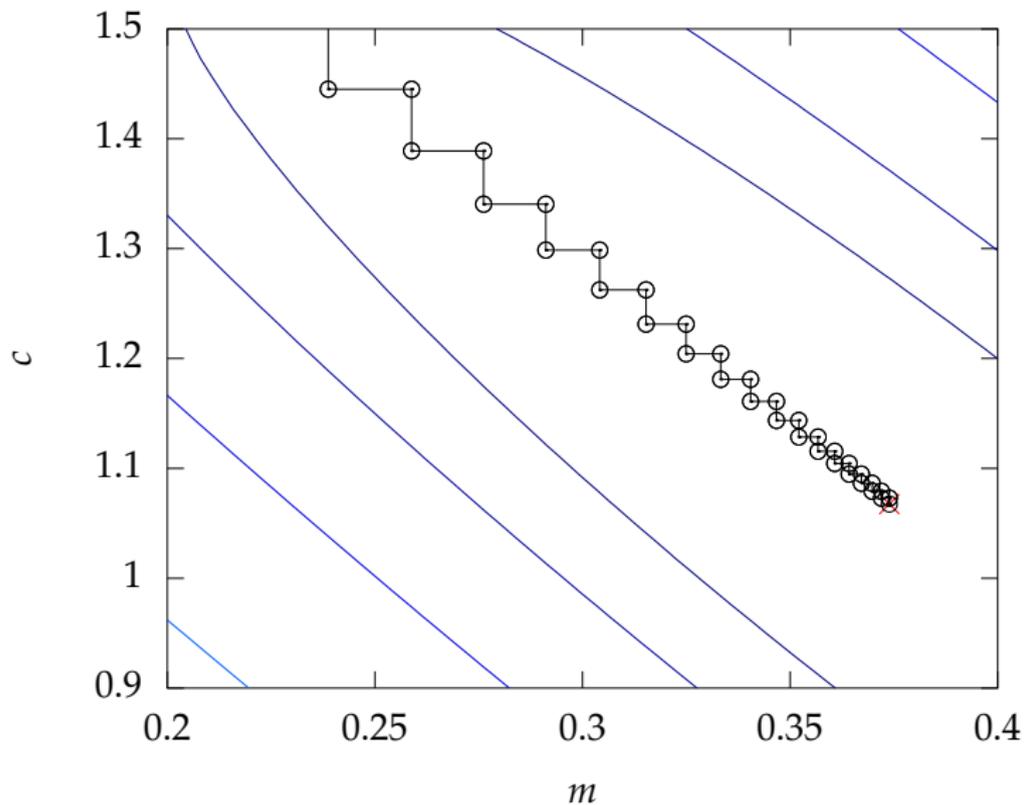
Coordinate Descent

Iteration 10



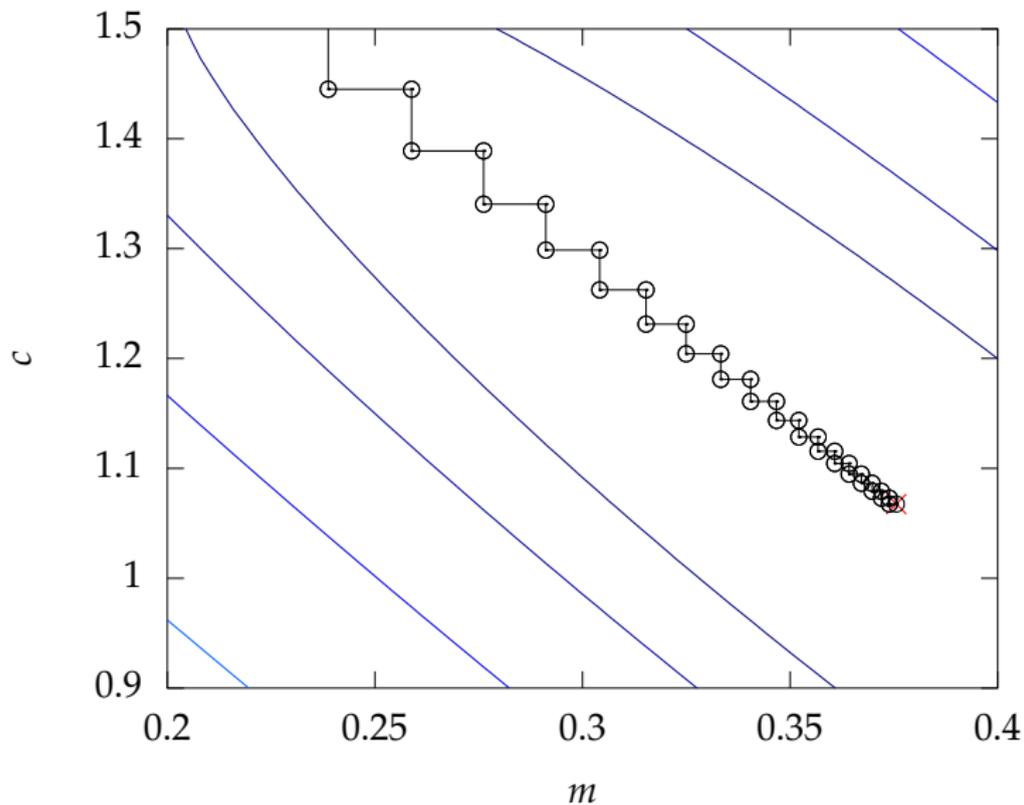
Coordinate Descent

Iteration 10



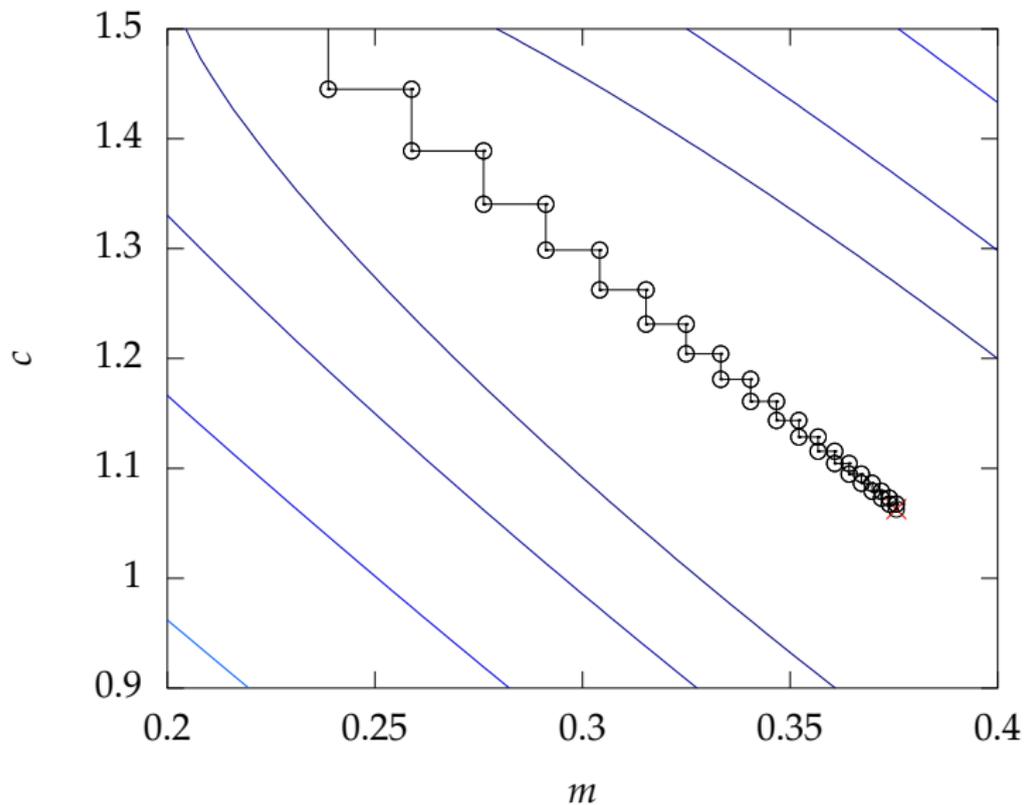
Coordinate Descent

Iteration 10



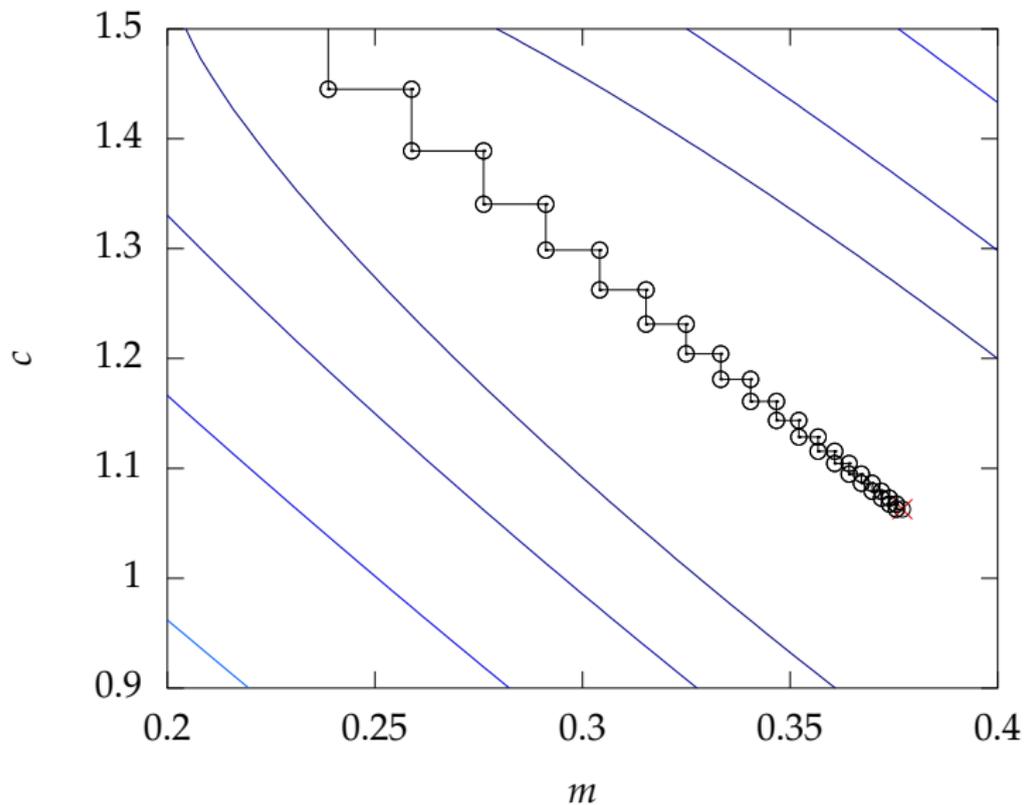
Coordinate Descent

Iteration 10



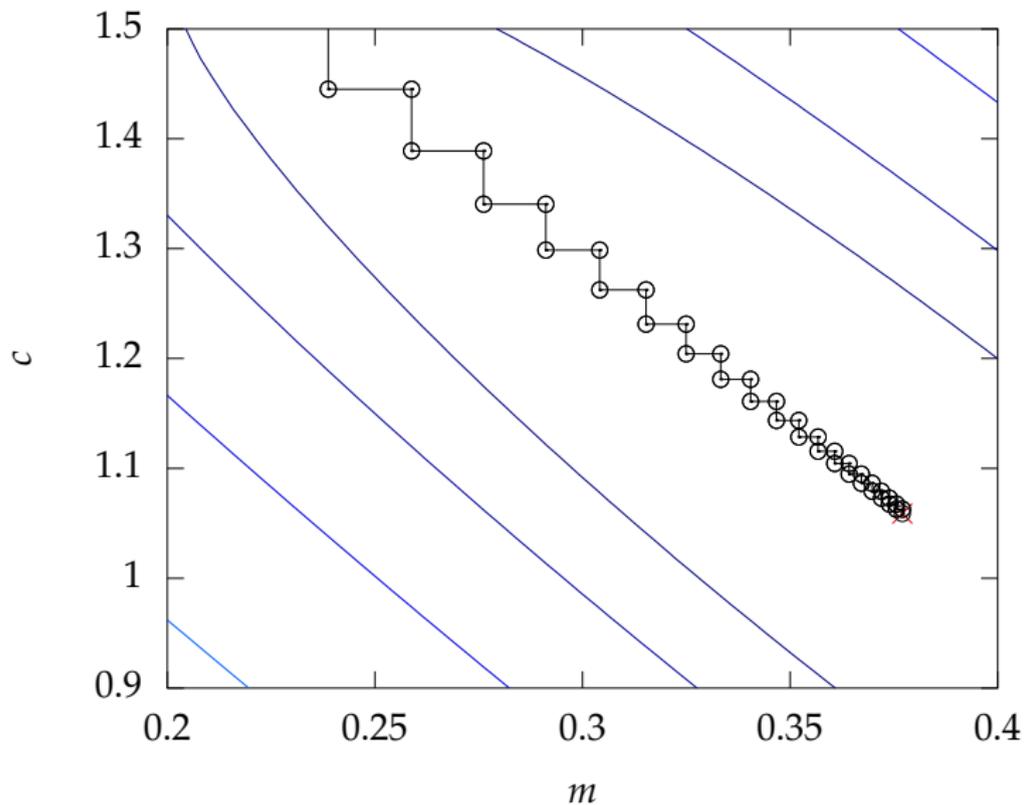
Coordinate Descent

Iteration 20



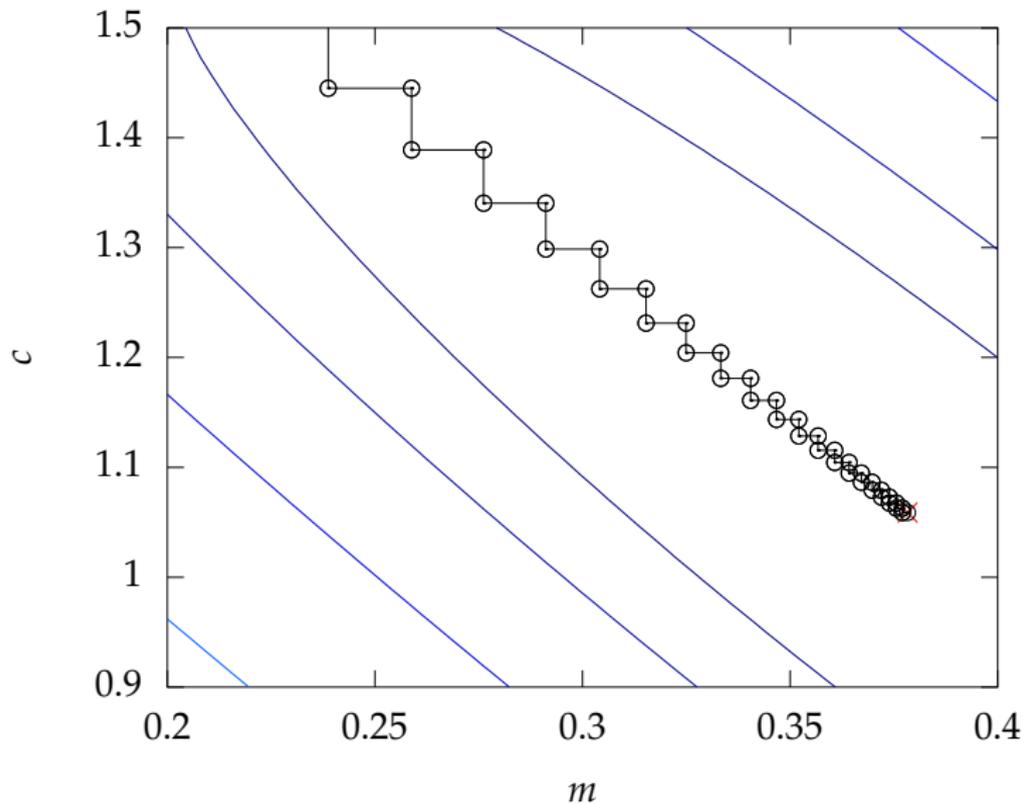
Coordinate Descent

Iteration 20



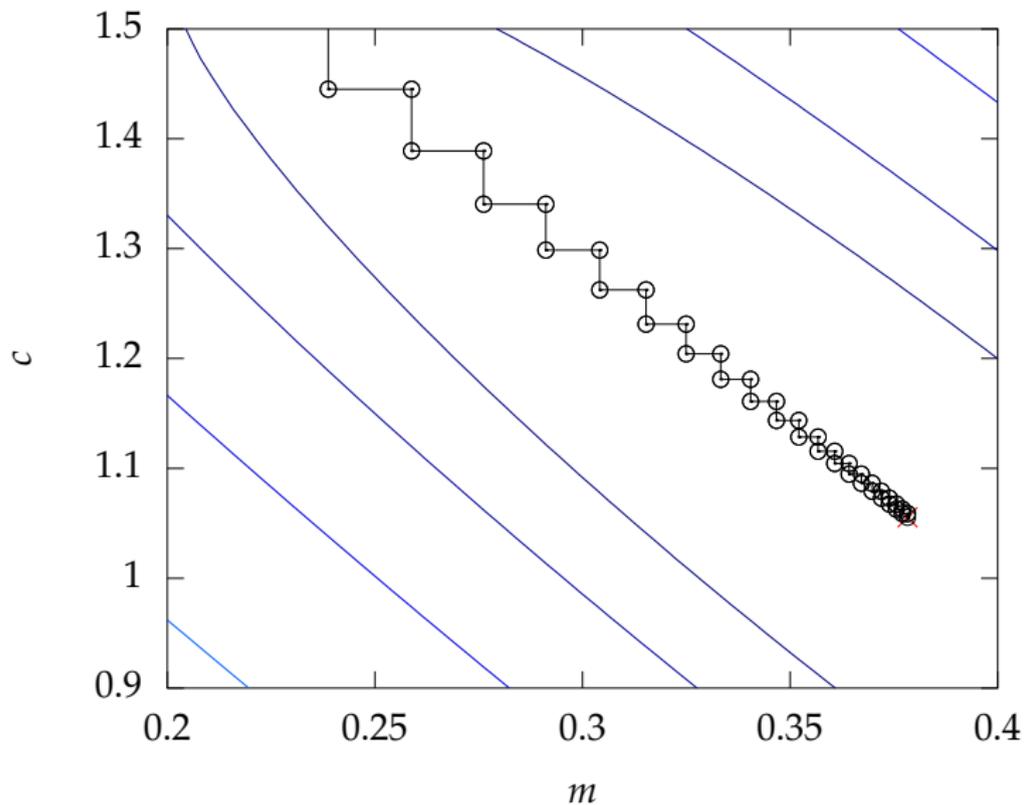
Coordinate Descent

Iteration 20



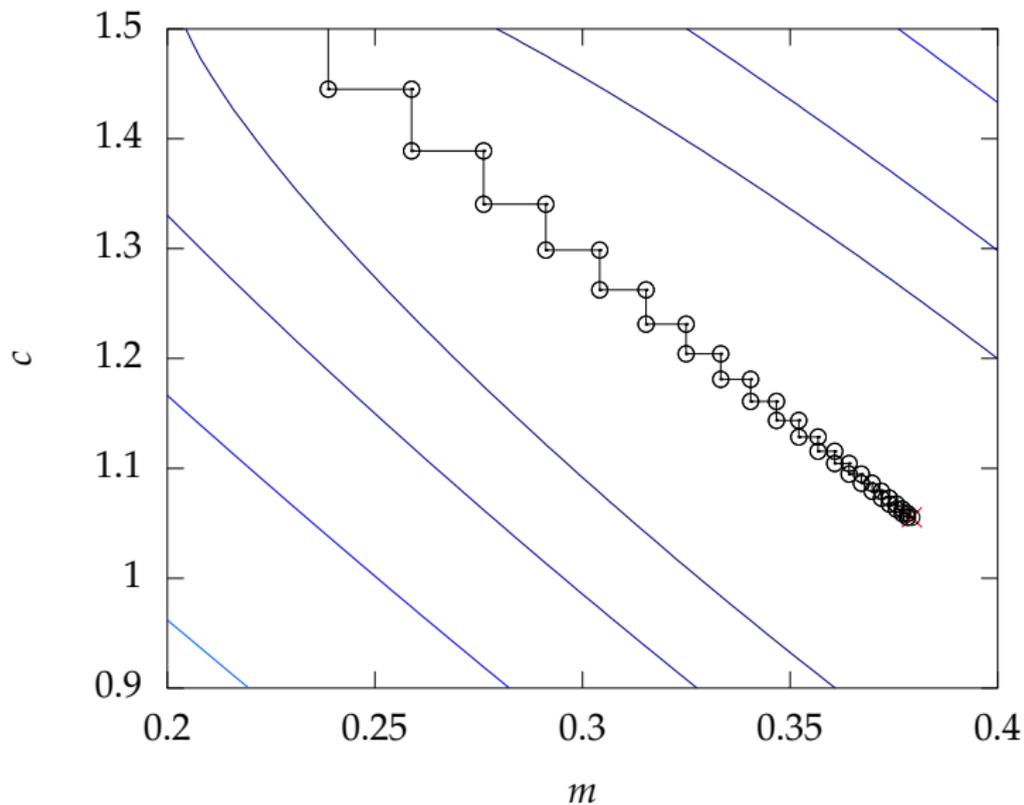
Coordinate Descent

Iteration 20



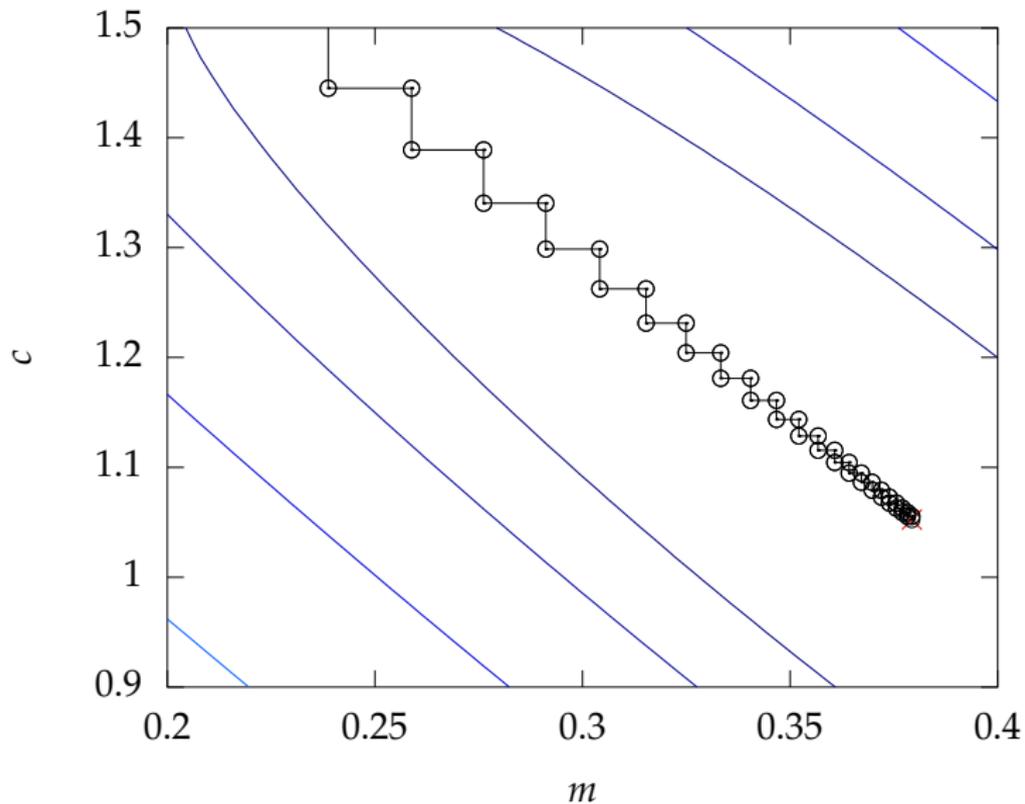
Coordinate Descent

Iteration 20



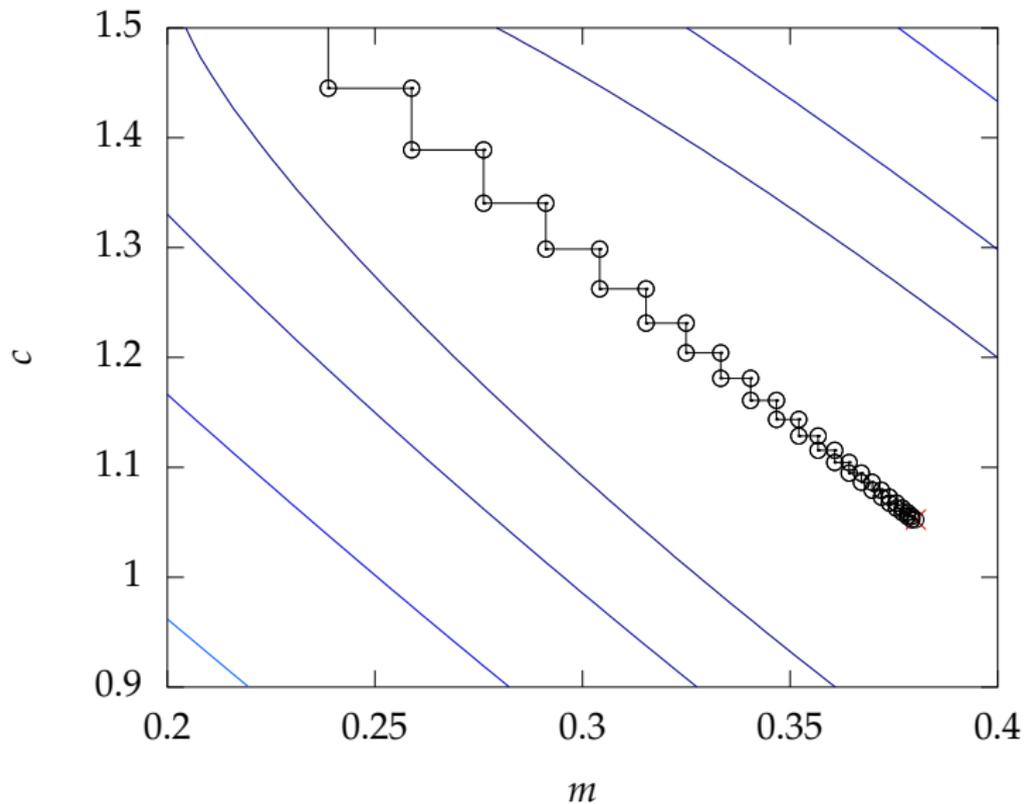
Coordinate Descent

Iteration 20



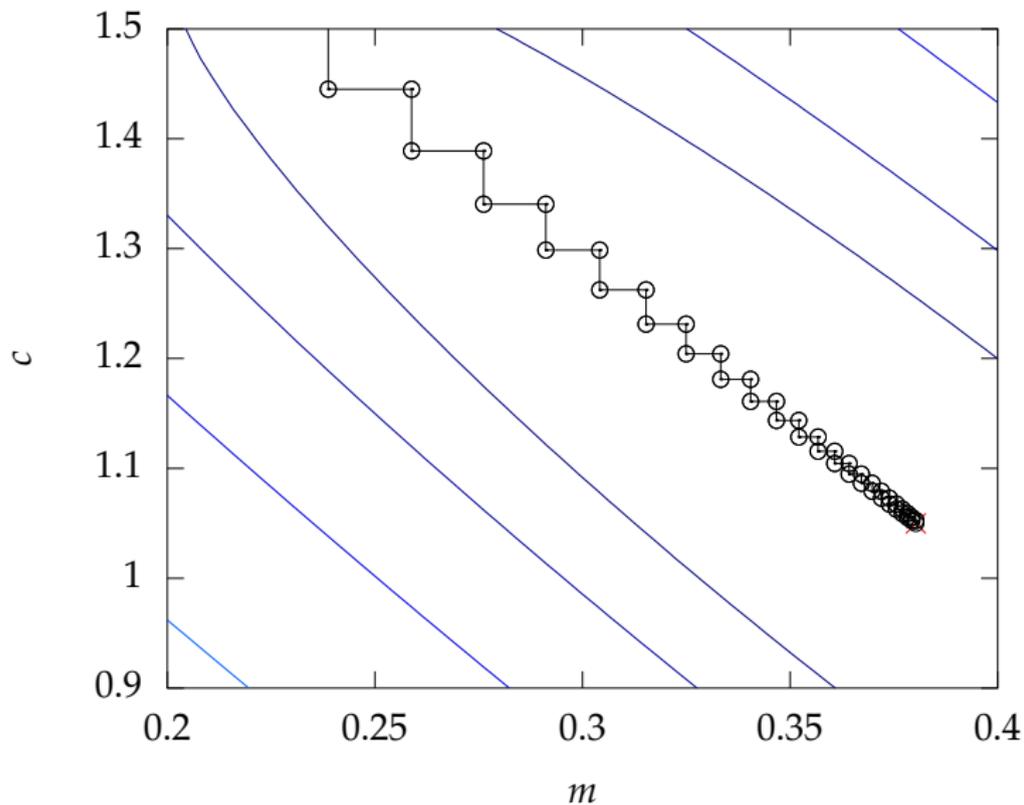
Coordinate Descent

Iteration 20



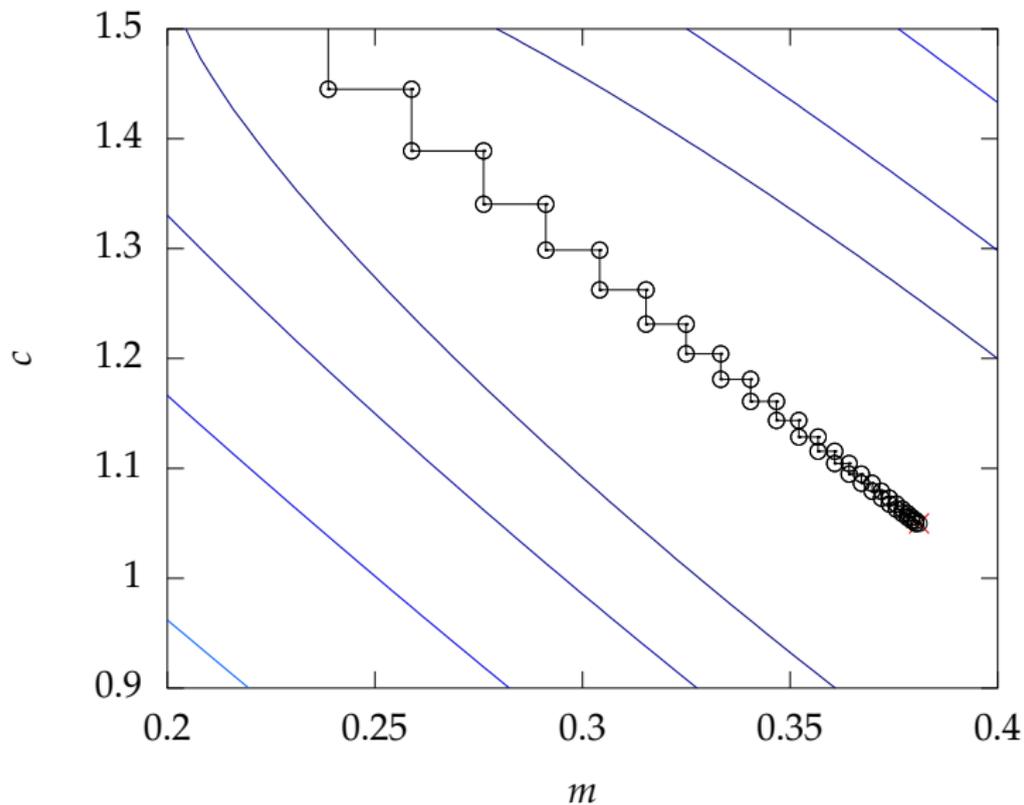
Coordinate Descent

Iteration 20



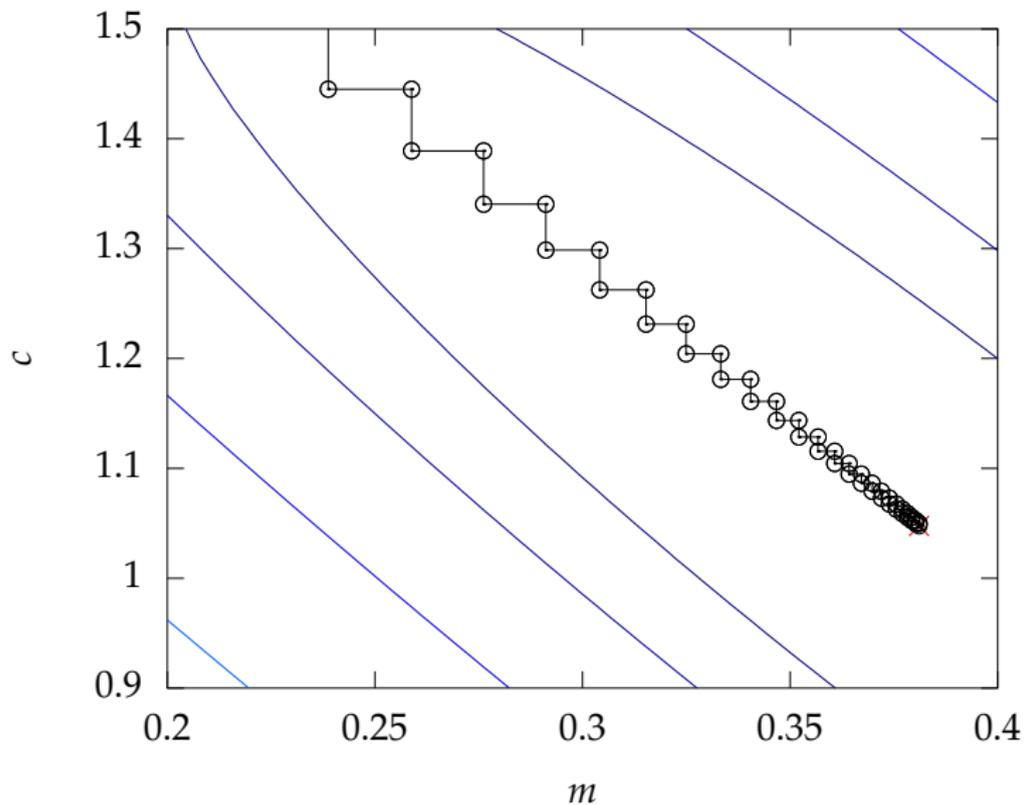
Coordinate Descent

Iteration 20



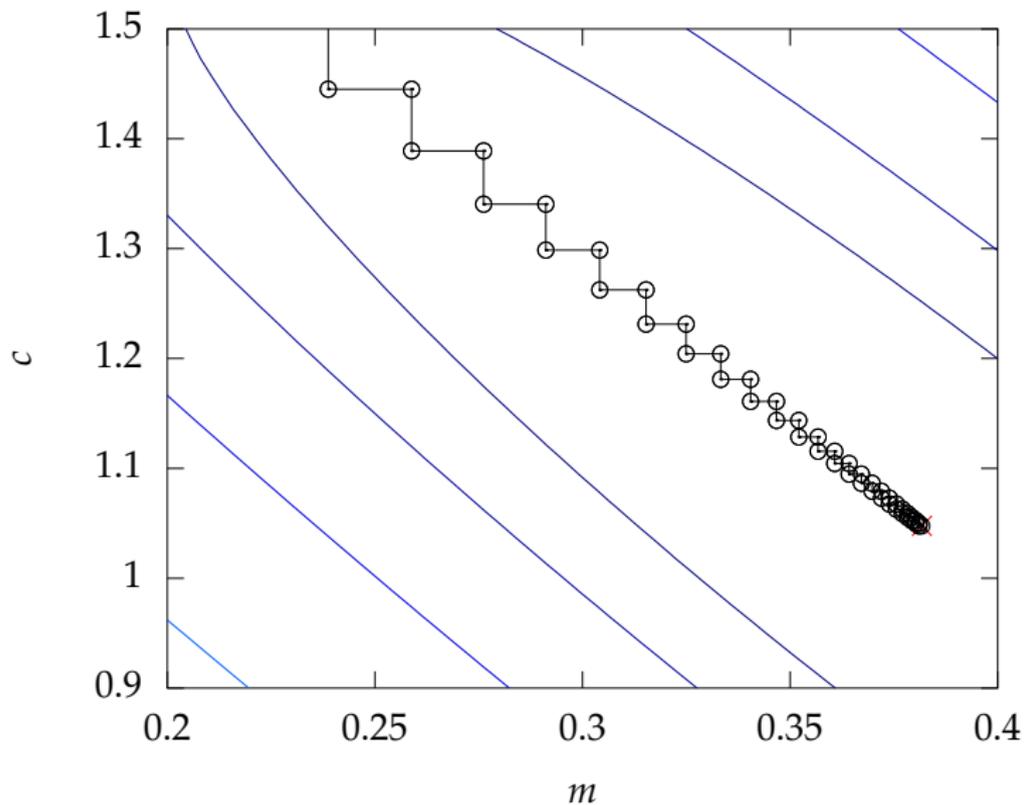
Coordinate Descent

Iteration 20



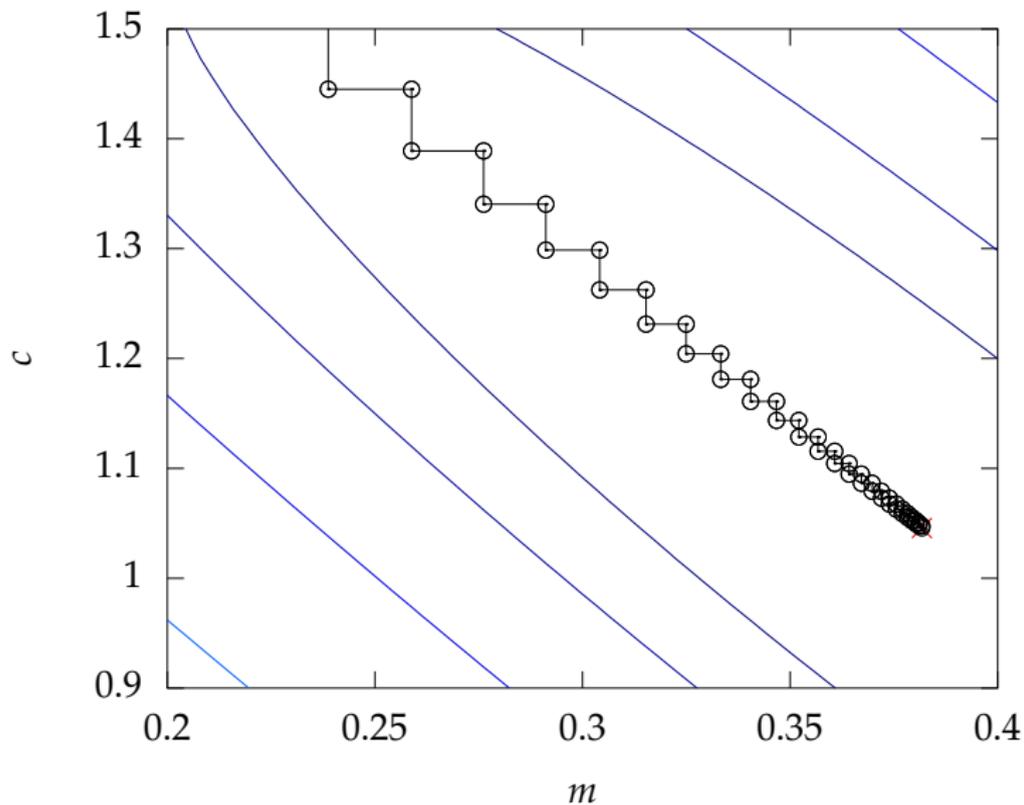
Coordinate Descent

Iteration 20



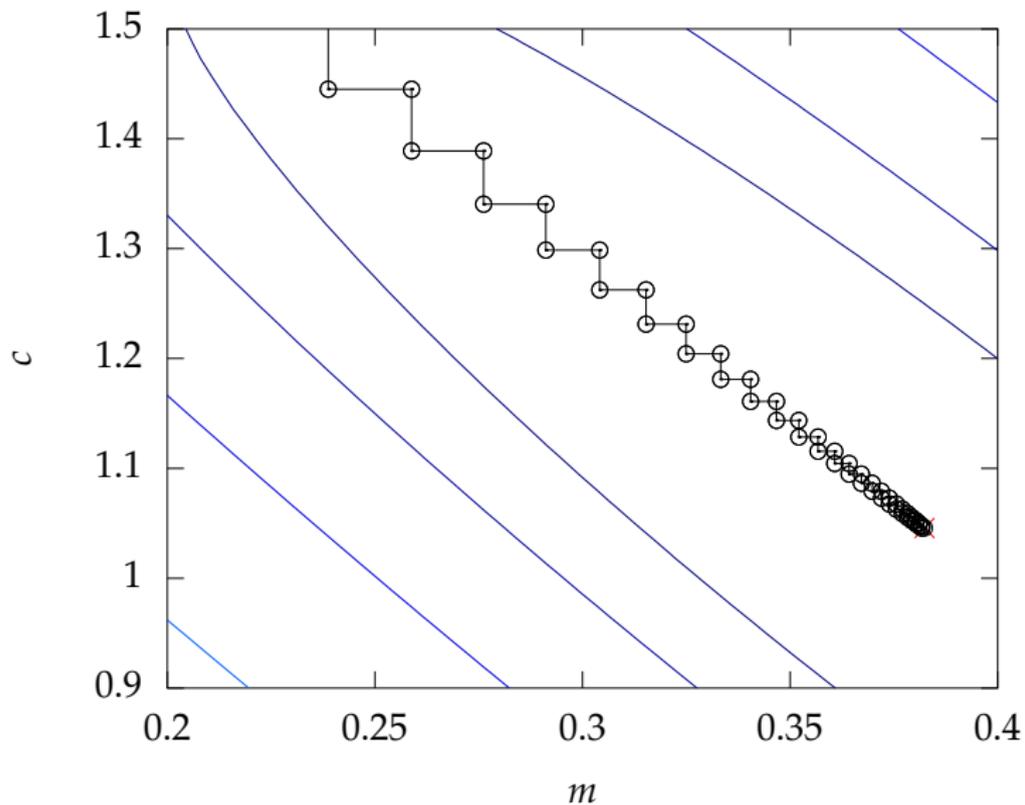
Coordinate Descent

Iteration 20



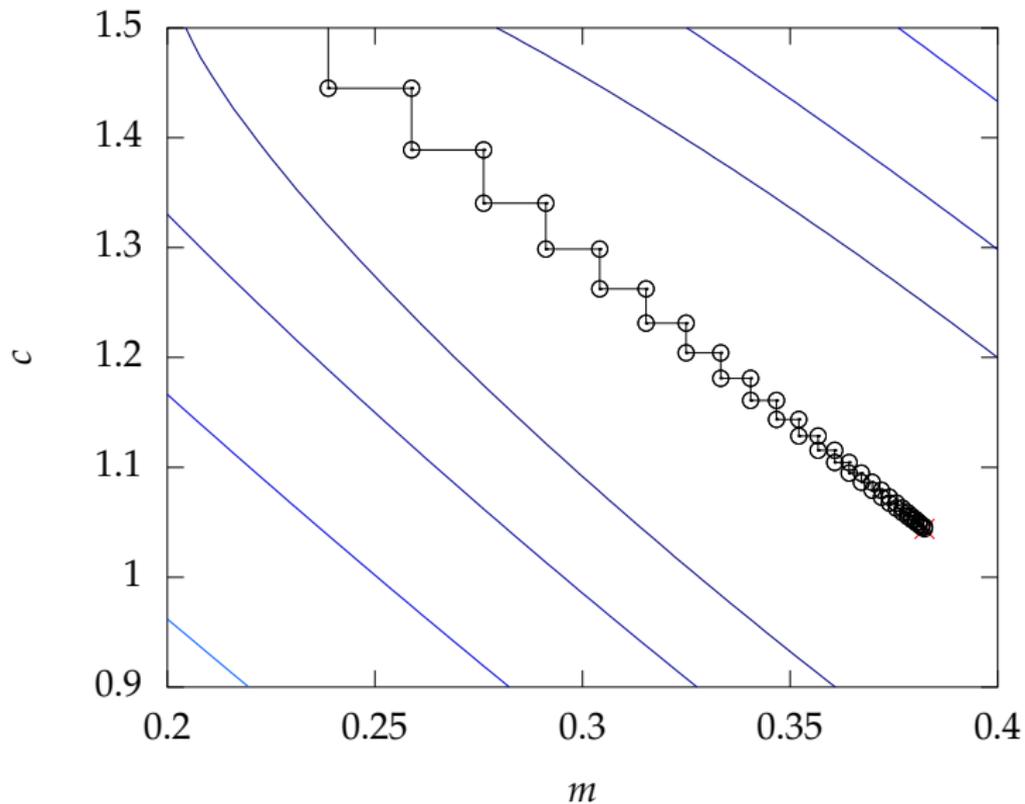
Coordinate Descent

Iteration 20



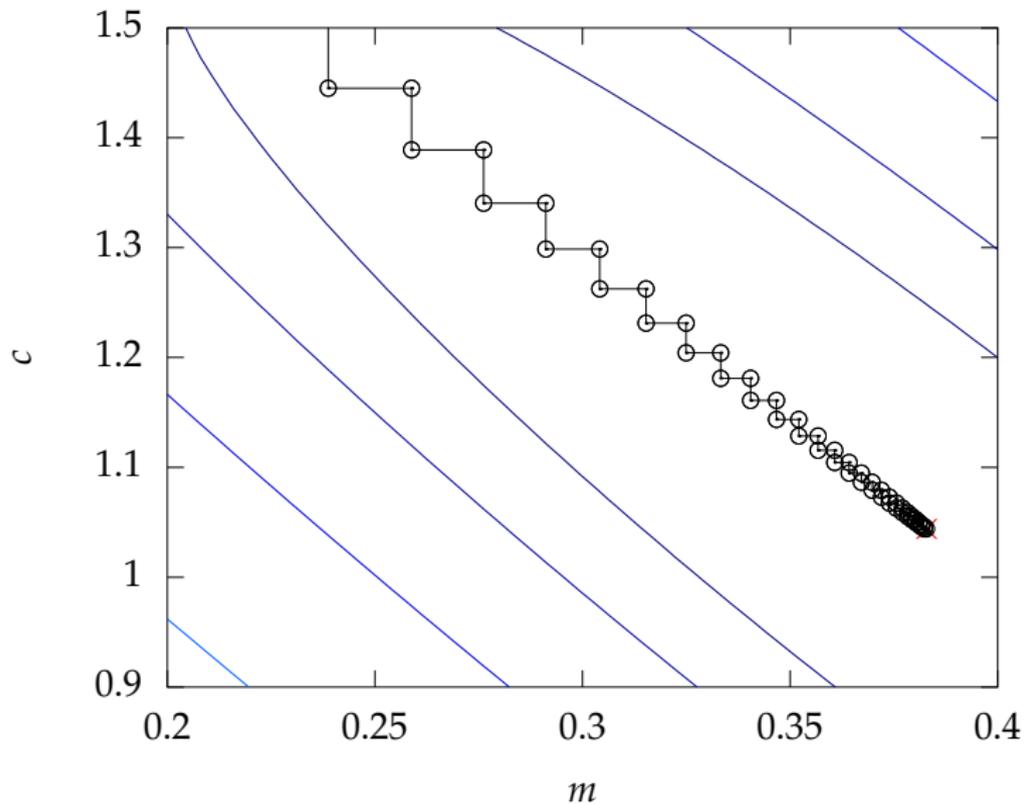
Coordinate Descent

Iteration 20



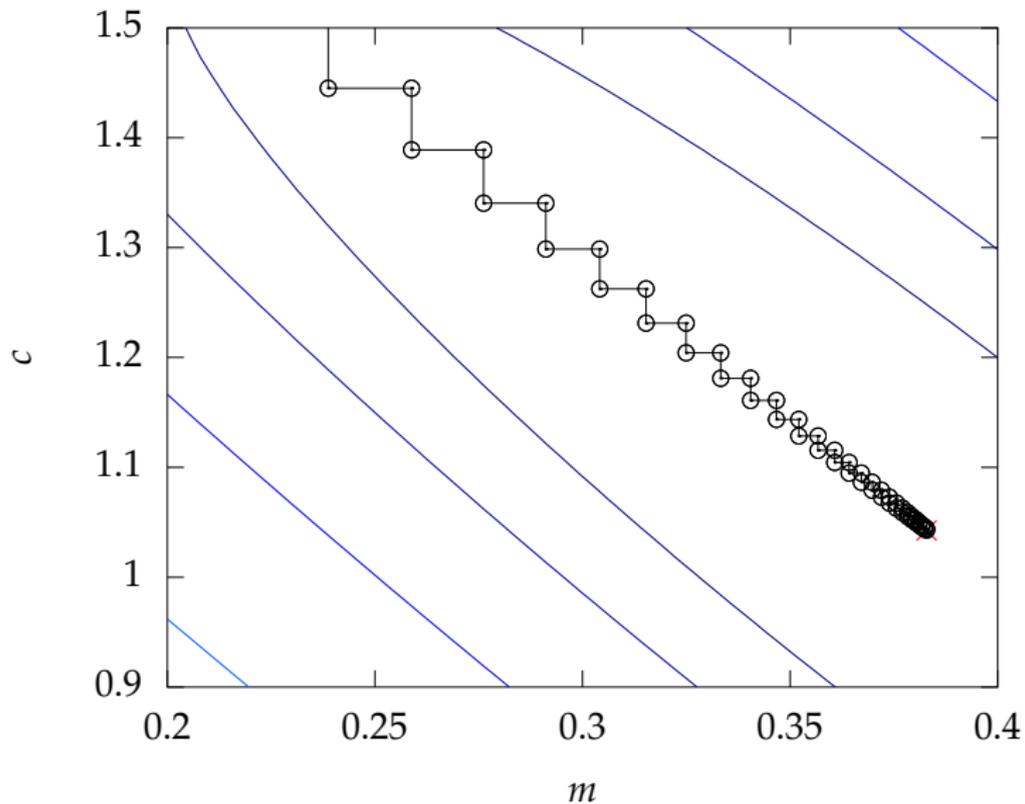
Coordinate Descent

Iteration 20



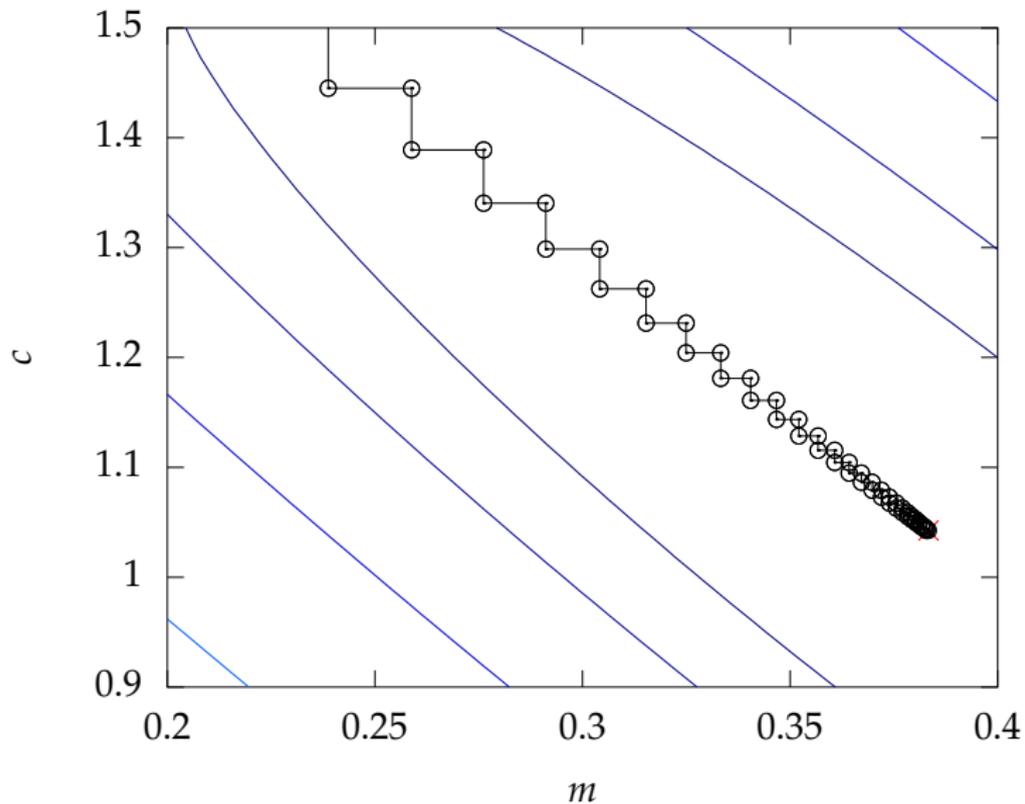
Coordinate Descent

Iteration 20



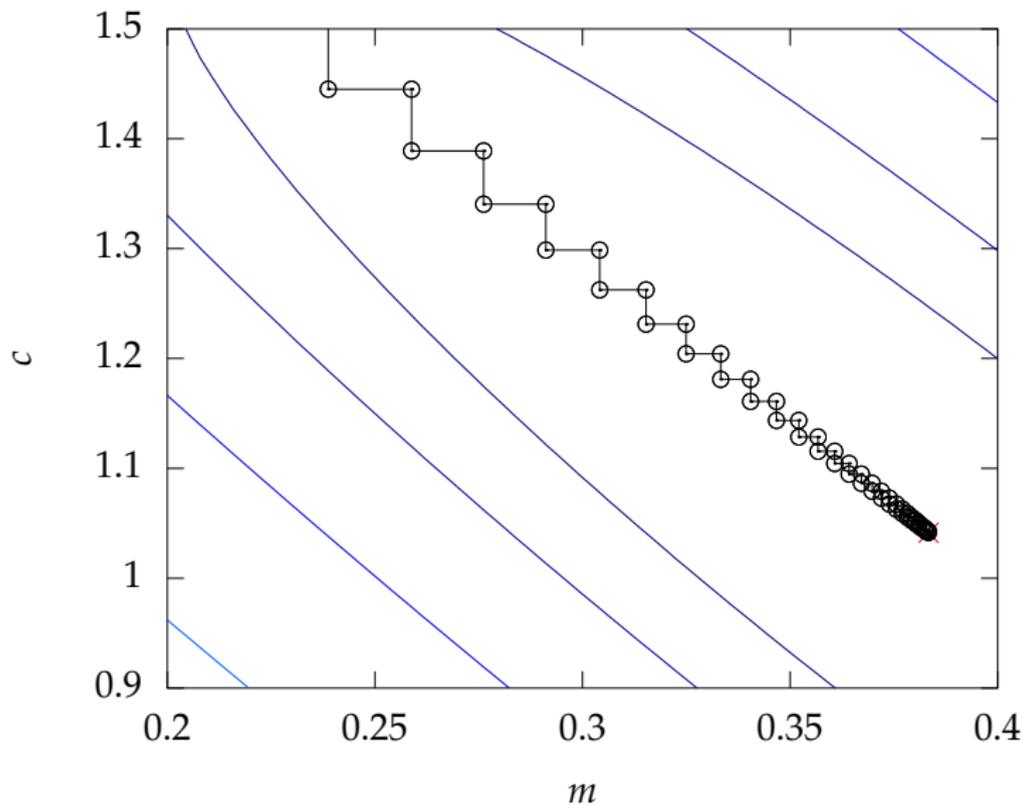
Coordinate Descent

Iteration 20



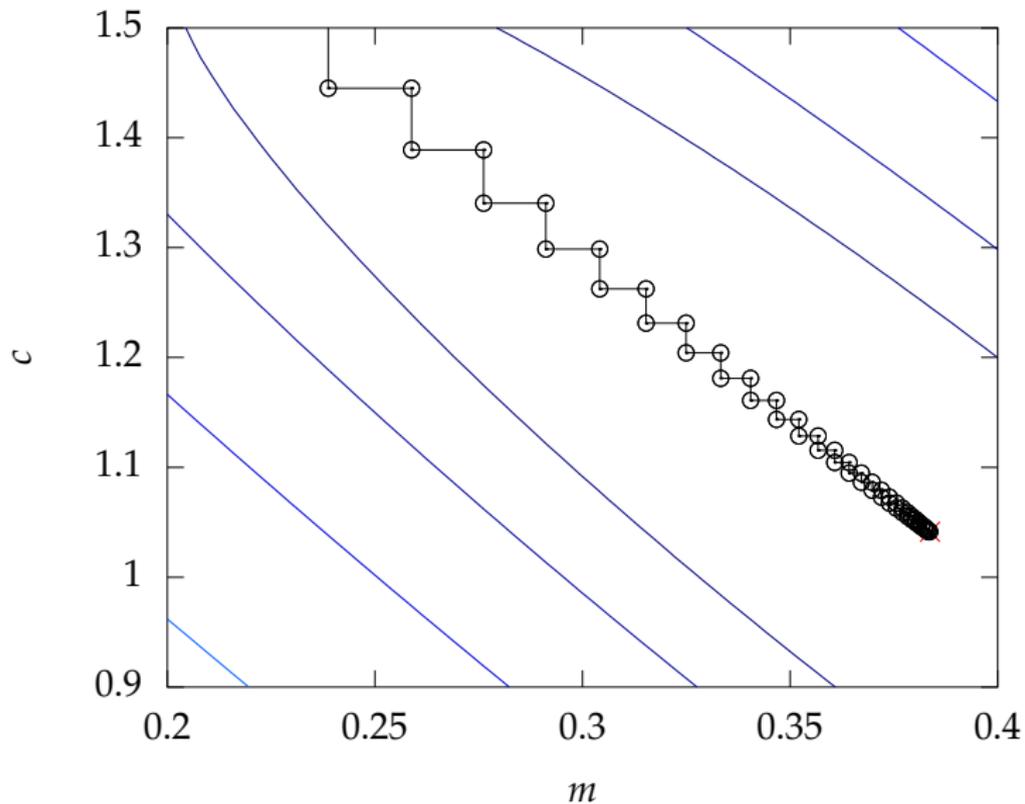
Coordinate Descent

Iteration 20



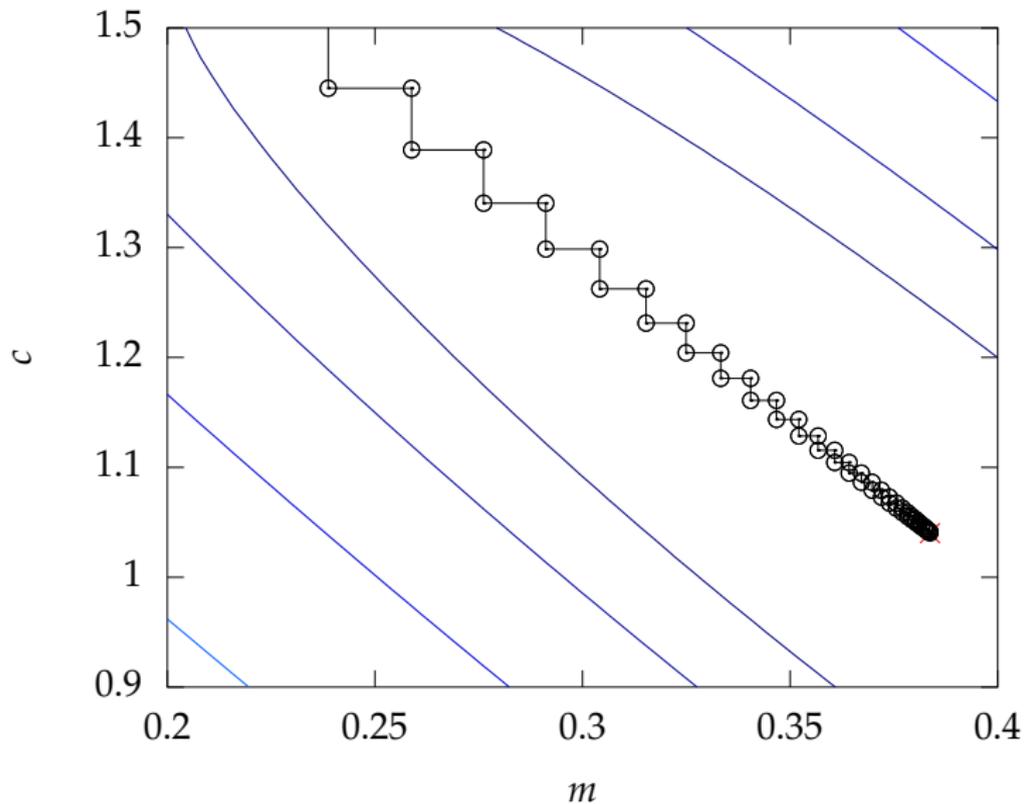
Coordinate Descent

Iteration 20



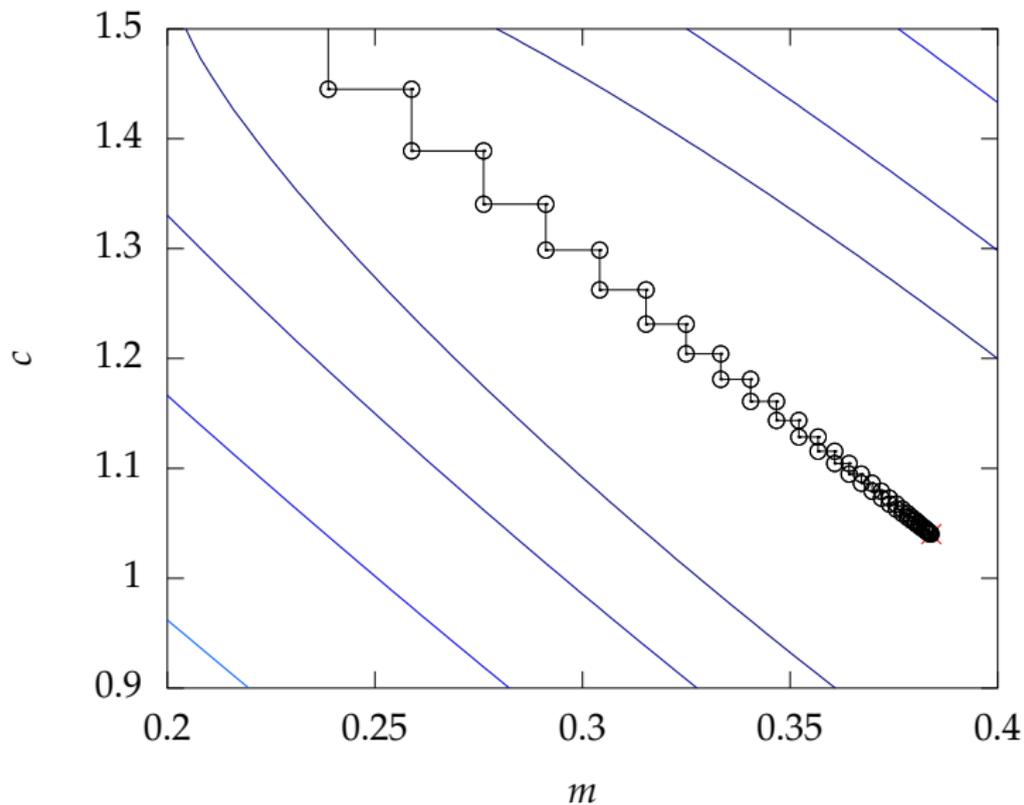
Coordinate Descent

Iteration 20



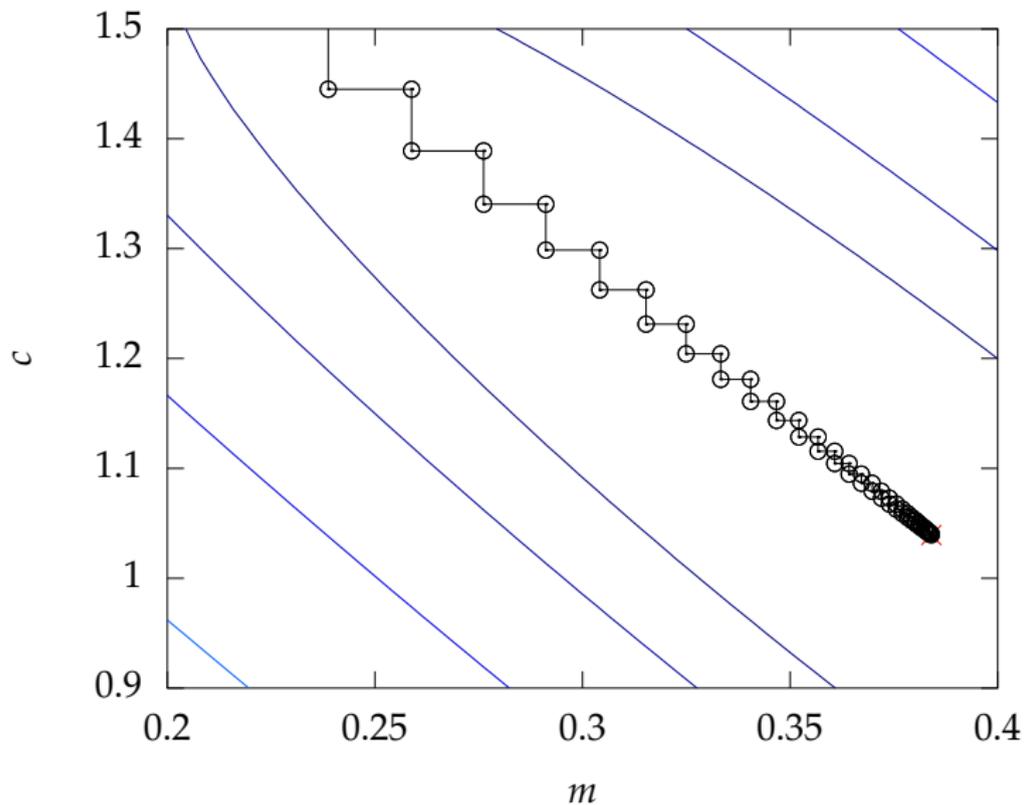
Coordinate Descent

Iteration 30



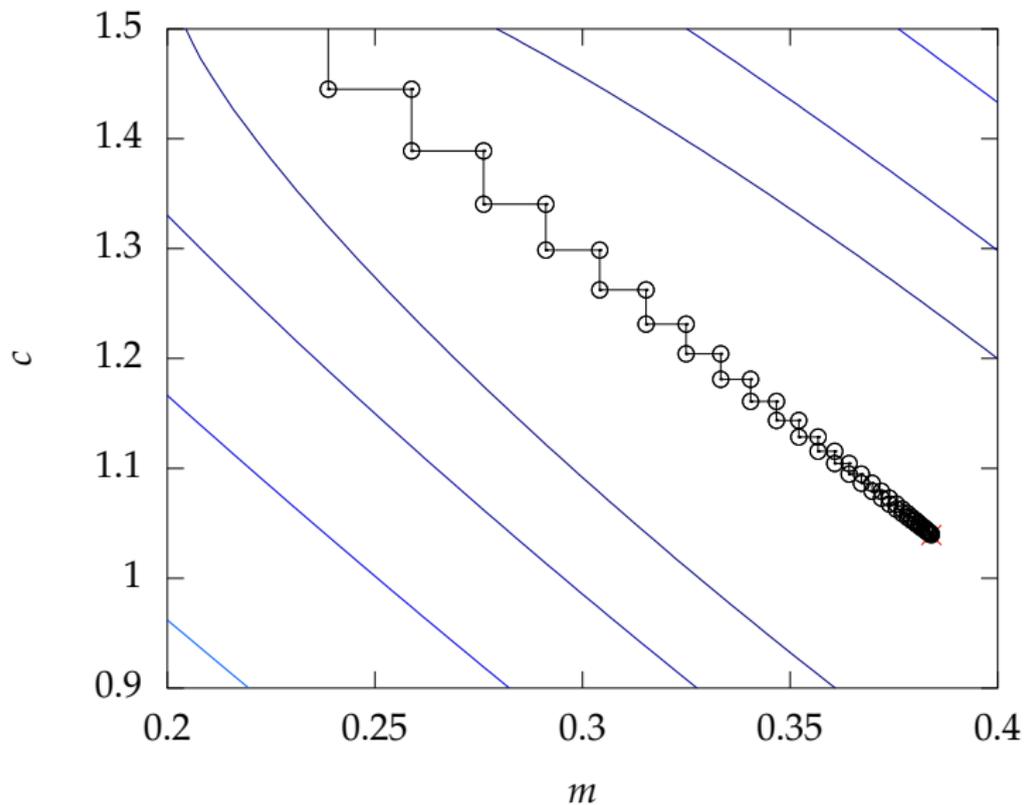
Coordinate Descent

Iteration 30



Coordinate Descent

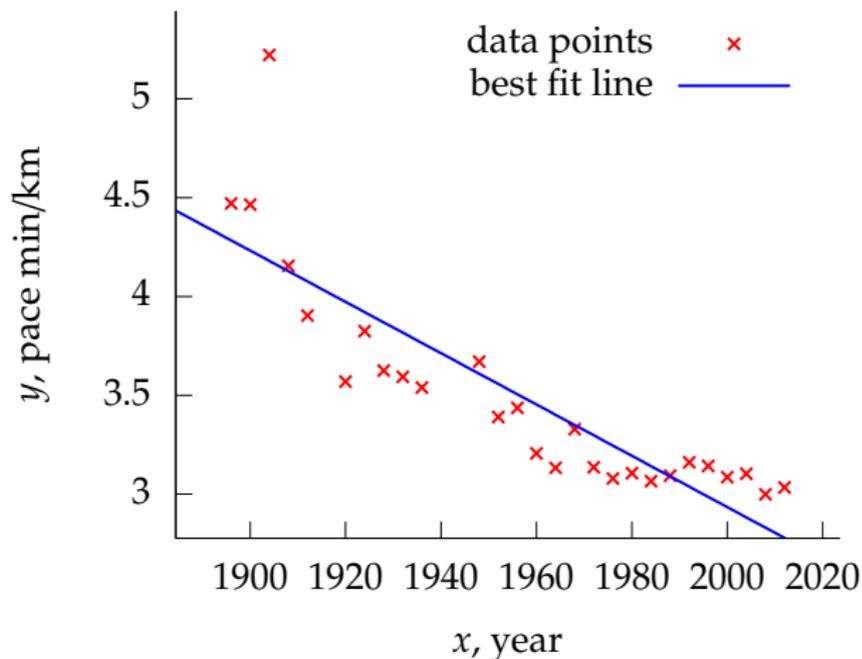
Iteration 30



Important Concepts Not Covered

- ▶ Optimization methods.
 - ▶ Second order methods, conjugate gradient, quasi-Newton and Newton.
 - ▶ Effective heuristics such as momentum.
- ▶ Local vs global solutions.

Linear Function



Linear regression for Male Olympics Marathon Gold Medal times.

Reading

- ▶ Section 1.1-1.2 of Rogers and Girolami for fitting linear models.
- ▶ Section 1.2.5 of Bishop up to equation 1.65.

Multi-dimensional Inputs

- ▶ Multivariate functions involve more than one input.
- ▶ Height might be a function of weight and gender.
- ▶ There could be other contributory factors.
- ▶ Place these factors in a feature vector \mathbf{x}_i .
- ▶ Linear function is now defined as

$$f(\mathbf{x}_i) = \sum_{j=1}^q w_j x_{i,j} + c$$

Vector Notation

mo

- ▶ Write in vector notation,

$$f(\mathbf{x}_i) = \mathbf{w}^\top \mathbf{x}_i + c$$

- ▶ Can absorb c into \mathbf{w} by assuming extra input x_0 which is always 1.

$$f(\mathbf{x}_i) = \mathbf{w}^\top \mathbf{x}_i$$

Log Likelihood for Multivariate Regression

- ▶ The likelihood of a single data point is

$$p(y_i|x_i) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - \mathbf{w}^\top \mathbf{x}_i)^2}{2\sigma^2}\right).$$

- ▶ Leading to a log likelihood for the data set of

$$L(\mathbf{w}, \sigma^2) = -\frac{n}{2} \log \sigma^2 - \frac{n}{2} \log 2\pi - \frac{\sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2}{2\sigma^2}.$$

- ▶ And a corresponding error function of

$$E(\mathbf{w}, \sigma^2) = \frac{n}{2} \log \sigma^2 + \frac{\sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2}{2\sigma^2}.$$

Expand the Brackets

$$\begin{aligned} E(\mathbf{w}, \sigma^2) &= \frac{n}{2} \log \sigma^2 + \frac{1}{2\sigma^2} \sum_{i=1}^n y_i^2 - \frac{1}{\sigma^2} \sum_{i=1}^n y_i \mathbf{w}^\top \mathbf{x}_i \\ &\quad + \frac{1}{2\sigma^2} \sum_{i=1}^n \mathbf{w}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{w} + \text{const.} \\ &= \frac{n}{2} \log \sigma^2 + \frac{1}{2\sigma^2} \sum_{i=1}^n y_i^2 - \frac{1}{\sigma^2} \mathbf{w}^\top \sum_{i=1}^n \mathbf{x}_i y_i \\ &\quad + \frac{1}{2\sigma^2} \mathbf{w}^\top \left[\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right] \mathbf{w} + \text{const.} \end{aligned}$$

Multivariate Derivatives

- ▶ We will need some multivariate calculus.
- ▶ For now some simple multivariate differentiation:

$$\frac{d\mathbf{a}^\top \mathbf{w}}{d\mathbf{w}} = \mathbf{a}$$

and

$$\frac{d\mathbf{w}^\top \mathbf{A} \mathbf{w}}{d\mathbf{w}} = (\mathbf{A} + \mathbf{A}^\top) \mathbf{w}$$

or if \mathbf{A} is symmetric (*i.e.* $\mathbf{A} = \mathbf{A}^\top$)

$$\frac{d\mathbf{w}^\top \mathbf{A} \mathbf{w}}{d\mathbf{w}} = 2\mathbf{A} \mathbf{w}.$$

Differentiate

Differentiating with respect to the vector \mathbf{w} we obtain

$$\frac{\partial L(\mathbf{w}, \beta)}{\partial \mathbf{w}} = \beta \sum_{i=1}^n \mathbf{x}_i y_i - \beta \left[\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right] \mathbf{w}$$

Leading to

$$\mathbf{w}^* = \left[\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right]^{-1} \sum_{i=1}^n \mathbf{x}_i y_i,$$

Rewrite in matrix notation:

$$\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top = \mathbf{X}^\top \mathbf{X}$$

$$\sum_{i=1}^n \mathbf{x}_i y_i = \mathbf{X}^\top \mathbf{y}$$

Update Equations

- ▶ Update for \mathbf{w}^* .

$$\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

- ▶ The equation for σ^{2*} may also be found

$$\sigma^{2*} = \frac{\sum_{i=1}^n (y_i - \mathbf{w}^{*\top} \mathbf{x}_i)^2}{n}.$$

- ▶ Section 1.3 of Rogers and Girolami for Matrix & Vector Review.

References I

- C. M. Bishop. *Pattern Recognition and Machine Learning*. Springer-Verlag, 2006. [\[Google Books\]](#) .
- P. S. Laplace. *Essai philosophique sur les probabilités*. Courcier, Paris, 2nd edition, 1814. Sixth edition of 1840 translated and reprinted (1951) as *A Philosophical Essay on Probabilities*, New York: Dover; fifth edition of 1825 reprinted 1986 with notes by Bernard Bru, Paris: Christian Bourgeois Éditeur, translated by Andrew Dale (1995) as *Philosophical Essay on Probabilities*, New York:Springer-Verlag.
- S. Rogers and M. Girolami. *A First Course in Machine Learning*. CRC Press, 2011. [\[Google Books\]](#) .