

## 7. BLACK-LITTERMAN OPTIMIZATION

**7.1. Statistics Review.** At a high level, forecasting may be viewed as a three-step procedure:

- (1) Build a model for the quantity you want to forecast, which connects that quantity to some parameters  $\theta$ .
- (2) Learn all you can about the parameter  $\theta$  from the data you have (perhaps using additional models which connect  $\theta$  to other quantities you can observe)
- (3) Use the model(s), and what you have learned about the parameter, to calculate the expected value of the quantity you want to forecast.

Here “what you have learned about the parameter” is encoded as a distribution on parameter space. This distribution is called the *posterior*.

For any probability measure  $p$  and any events  $A, B$ ,

$$\begin{aligned} p(A|B)p(B) &= p(A \cap B) = p(B \cap A) = p(B|A)p(A) \\ p(A|B) &= p(B|A) \frac{p(A)}{p(B)}. \end{aligned} \quad (7.1)$$

This result, known as *Bayes' theorem*, appears trivial but historically was a major conceptual step because it allows the “inversion” of probabilities.

Eq. (7.1) can be applied to any events  $A$  and  $B$ , but is most usefully applied when  $A$  is the event of observing a certain data set (usually the one given to us by the world), and  $B$  is the event of parameters taking on certain values.

The following is standard terminology in the statistics literature (Robert, 2007). A *Bayesian statistical model* consists of:

- (1) a vector-valued random variable  $\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^d$  distributed according to  $f(\mathbf{x}|\theta)$ , where realizations of  $\mathbf{x}$  have been observed and only the parameter  $\theta$  (which belongs to a real vector space  $\Theta \subseteq \mathbb{R}^\ell$ ) is unknown, and
- (2) a prior density  $\pi(\theta)$  on  $\Theta$ .

The function  $f(\mathbf{x}|\theta)$  is called the *likelihood* and, after conditioning on  $\theta$ , forms a density on the *data space*  $\mathcal{X} \subseteq \mathbb{R}^d$ . Applying Bayes' theorem to the likelihood function gives the *posterior*

$$p(\theta|\mathbf{x}) = \frac{f(\mathbf{x}|\theta)\pi(\theta)}{p(\mathbf{x})} = \frac{\text{likelihood} \times \text{prior}}{\text{normalization}}$$

The posterior is best viewed as the density on  $\Theta$  proportional to  $f(\mathbf{x}|\theta)\pi(\theta)$ , and the normalization factor  $p(\mathbf{x})$  drops out of certain calculations. In Bayesian statistics, all statistical inference is based on the posterior.

*Example 7.1* (Bayes (1763)). A ball  $A$  is rolled along the unit interval  $[0, 1]$ , with a uniform probability of stopping at any point. It stops at a point  $\theta$  between 0 and 1, and is not moved subsequently. A second ball  $B$  is then rolled  $n$  times under the same assumptions. Let  $x$  denote the number of times the ball  $B$  stopped before passing  $A$ . Given  $x$ , what inference can we make on  $\theta$ ?

Once  $\theta$  is known, stopping to the left versus the right of  $\theta$  is just like flipping an unfair coin with  $\text{prob}(\text{heads}) = \theta$ . Experiments like flipping an unfair coin are known in mathematics as *Bernoulli trials*.

Formally, a *Bernoulli trial* is any random experiment with only two outcomes, usually called “success” and “failure” and conveniently denoted by 1 or 0. A *Bernoulli process* is a finite or infinite sequence of independent Bernoulli trials with the same probability of success, ie. a sequence of independent random variables  $X_1, X_2, X_3, \dots$  with  $X_i \in \{0, 1\}$ , and there exists some  $\theta \in [0, 1]$  such that  $P(X_i = 1) = \theta$  for all  $i \in \mathbb{N}$ . The number of successes in  $n$  independent Bernoulli trials is a random variable whose distribution is called the *Binomial distribution* and denoted  $B(n, \theta)$ .

Therefore, in Example 7.1,  $x \sim B(n, \theta)$ . Bayes’ theorem tells us

$$p(\theta | x) \propto p(x | \theta)p(\theta). \quad (7.2)$$

By assumption  $\theta$  is uniform on  $[0, 1]$ , so  $p(\theta) = 1$ . Hence (7.2) is proportional to the binomial density  $p(x | \theta) \propto \theta^x (1 - \theta)^{n-x}$  considered as a function of  $\theta$ , which is known as the *beta distribution*  $\text{Beta}(x + 1, n - x + 1)$ .

This example illustrates that the prior isn’t necessarily something connected with opinions or views. Sometimes it is dictated by physics. In this example, the uniform prior is determined by the mechanics of the process which occurred first (the rolling of the first ball).

**7.2. Prediction.** In particular, given any quantity of interest  $r$  which is related to  $\theta$  by some other density  $p(r | \theta)$ , with the posterior in hand we can form a prediction of  $r$ :

$$\mathbb{E}[r | x] = \int r \underbrace{\int p(r | \theta) p(\theta | x) d\theta}_{\text{posterior predictive density}} dr$$

$p(r, \theta | x)$   
 $\times \quad \nwarrow \quad \nearrow \quad Y$   
 $\theta$

Let’s consider prediction from the unfair-coin model. What is the probability that the next flip, say  $y$ , is heads (denoted by 1 for convenience).

$$p(y = 1 | x) = \int p(y = 1 | \theta, x) p(\theta | x) d\theta = \frac{x + 1}{n + 2}.$$

This is known as “Laplace’s law of succession.” Note that

$$p(y = 1 | \theta, x) = \theta,$$

so our forecast is just the posterior mean. The mean of  $\text{Beta}(\alpha, \beta)$  is  $\alpha/(\alpha + \beta)$ .

**7.3. Conjugate priors.** Now suppose that, instead of the uniform prior implied by Bayes' example, we take  $\text{Beta}(\alpha, \beta)$  as our prior. In other words, we take as prior the following:

$$p(\theta | \alpha, \beta) \propto \theta^{\alpha-1}(1-\theta)^{\beta-1}$$

The posterior is then

$$\begin{aligned} p(\theta | x) &\propto p(x | \theta)p(\theta) \\ &\propto \theta^x(1-\theta)^{n-x}\theta^{\alpha-1}(1-\theta)^{\beta-1} \\ &\sim \text{Beta}(\alpha + x, \beta + n - x). \end{aligned}$$

也就是说 likelihood 部分是 binomial, prior 和 posterior 可以都是 beta distribution, 这叫做 conjugate.

Thus the posterior is of the same form as the prior, but with different parameters.

**The beta distribution is said to be a conjugate prior for the binomial model.**

We also say the prior has been “updated” by the data,  $x$  in this case. If more data come in, the posterior is suitable for use as a prior in the next round of updating. The order in which i.i.d. observations are collected does not matter, and updating the prior one observation at a time, or all observations together, does not matter. This kind of model lends itself well to inference with big data sets, since all of the relevant information about the full data set – no matter how large it may be – is contained in a small number of parameters.

**7.4. The Normal Distribution.** Old friends are the best kind. Specifically, consider an i.i.d. sample  $x = (x_1, \dots, x_n)$  drawn from a normal distribution  $N(\mu, \sigma^2)$ . It is convenient to write  $\tau = 1/\sigma^2$ . In this course we will generally always use  $\theta$  to denote the vector of parameters in the model. In this case  $\theta = (\mu, \tau)$ .

Consider the prior

$$\tau \sim \Gamma(\alpha, \beta), \quad \mu | \tau \sim N\left(\nu, \frac{1}{k\tau}\right) \text{ for } k > 0, \nu \in \mathbb{R}.$$

This prior distribution for  $\tau$  involves the Gamma distribution, which exists in several different parameterizations. The one we'll use has pdf given by

$$\beta^\alpha x^{\alpha-1} e^{-x\beta} / \Gamma(\alpha) \text{ for } x \geq 0 \text{ and } \alpha, \beta > 0.$$

In this case  $\alpha$  is called the **shape parameter**, while  $\beta$  is called the **inverse-scale or rate parameter**.

If  $\tau \sim \Gamma(\alpha, \beta)$  with these conventions, this means that  $\sigma^2 \sim \text{IG}(\alpha, \beta)$  follows the aptly-named inverse-gamma distribution. The inverse gamma distribution's pdf is defined over the support  $x > 0$  and is given by

$$f(x; \alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{-\alpha-1} \exp\left(-\frac{\beta}{x}\right).$$

Note that to get the density of  $\sigma^2$  from the density of  $\tau = 1/\sigma^2$  one would use the change-of-variables formula. Aside: This distribution is not only useful as a Bayesian prior; the first hitting time of a standard Wiener process has an inverse-gamma distribution with parameters  $\alpha = 1/2$  and  $\beta = x^2/2$  where  $x$  is the value to hit.

The full prior density is

$$\pi(\tau, \mu) = \pi(\tau)\pi(\mu | \tau),$$

which may be written as

$$\pi(\tau, \mu) = \frac{\beta^\alpha}{\Gamma(\alpha)} \tau^{\alpha-1} e^{-\beta\tau} \cdot \sqrt{\frac{k\tau}{2\pi}} \exp\left\{-\frac{k\tau}{2}(\mu - \nu)^2\right\}.$$

or more simply

$$\pi(\tau, \mu) \propto \tau^{\alpha-1/2} \exp\left[-\tau\left\{\beta + \frac{k}{2}(\mu - \nu)^2\right\}\right].$$

We have  $x_1, \dots, x_n$  independent, identically distributed from  $N(\mu, 1/\tau)$ , so the likelihood is

$$p(x | \mu, \tau) = \left(\frac{\tau}{2\pi}\right)^{n/2} \exp\left\{-\frac{\tau}{2} \sum_i (x_i - \mu)^2\right\}$$

Thus

$$\pi(\tau, \mu | x) \propto \tau^{\alpha+n/2-1/2} \exp\left[-\tau\left\{\beta + \frac{k}{2}(\mu - \nu)^2 + \frac{1}{2} \sum_i (x_i - \mu)^2\right\}\right].$$

Complete the square to see that

$$\begin{aligned} k(\mu - \nu)^2 + \sum_i (x_i - \mu)^2 &= (k+n) \left(\mu - \frac{k\nu + n\bar{x}}{k+n}\right)^2 \\ &\quad + \frac{nk}{n+k}(\bar{x} - \nu)^2 + \sum_i (x_i - \bar{x})^2. \end{aligned}$$

Hence the posterior satisfies

$$\pi(\tau, \mu | x) \propto \tau^{\alpha'-1/2} \exp\left[-\tau\left\{\beta' + \frac{k'}{2}(\mu - \nu')^2\right\}\right]$$

where

$$\begin{aligned} \alpha' &= \alpha + \frac{n}{2} \\ \beta' &= \beta + \frac{1}{2} \frac{nk}{n+k}(\bar{x} - \nu)^2 + \frac{1}{2} \sum_i (x_i - \bar{x})^2 \\ k' &= k + n \\ \nu' &= \frac{k\nu + n\bar{x}}{k+n}. \end{aligned}$$

收集关于  $\mu$  的一次项  
关于  $\mu$  的一次项  
和常数项  
 $a(\mu-b)^2 = a\mu^2 - 2ab\mu + ab^2$   
 $\mu^2$  的系数就是  $a$   
 $b$  可由  $\mu$  的系数  $-2a$  得到  
再按常数项.

原 Assumption 中  
 $\mu | \tau \sim N(\nu, \frac{1}{k\tau})$   
由于加入新数据  $\mu | \tau$  的原  
方差由  $\frac{\sigma^2}{k} \rightarrow \frac{\sigma^2}{n+k}$ , 变小  
原均值变为  $\nu$  与  $\bar{x}$  的 weighted  
average

$$\begin{aligned} &= \frac{k}{k+n} \nu + \frac{n}{k+n} \bar{x} \\ &\quad \downarrow \quad \quad \downarrow \\ &\text{Assumption} \quad \text{样本均值} \end{aligned}$$

Thus the posterior distribution is of the same parametric form as the prior (the above form of prior is a conjugate family), but with  $(\alpha, \beta, k, \nu)$  replaced by  $(\alpha', \beta', k', \nu')$ . To write code to update this model with new data, all you really need is to keep track of the sufficient statistics! Equivalently you could keep track of the parameters  $(\alpha, \beta, k, \nu)$ .

**7.5. Bayesian Multivariate Regression.** We now proceed with a fully Bayesian multivariate analysis of the likelihood

$$p(y | \theta) = (2\pi\sigma^2)^{-n/2} \exp \left[ -\frac{1}{2\sigma^2} \|y - X\beta\|^2 \right] \quad (7.3)$$

After completing the square, (7.3) can be written as follows:

$$\ell(\beta, \sigma^2 | y) = (2\pi\sigma^2)^{-n/2} \exp \left[ -\frac{1}{2\sigma^2} \left\{ s^2 + (\beta - \hat{\beta})'(X'X)(\beta - \hat{\beta}) \right\} \right] \quad (7.4)$$

Eq. (7.4) implies that  $T(y) = (s^2, \hat{\beta})$  forms a  $(k+1)$ -dimensional (minimal) sufficient statistic for the parameters  $\theta = (\sigma^2, \beta)$ .

In real-world problems with large data sets, we are concerned with “online” updating. In an online scenario, we are concerned with how we will update the sufficient statistics if a new row of  $X$  and a new response variable are published. The answer to this question is that we can update the  $k + k(k+1)/2$  numbers  $X'X$  and  $X'y$ , and use those to update the minimal sufficient statistics  $(\hat{\beta}, s^2)$ .

A conjugate prior can be constructed within the Gaussian-Inverse-Gamma family:

$$\text{prior: } \beta | \sigma^2 \sim N(\beta_0, \sigma^2 M^{-1}), \quad \sigma^2 \sim \text{IG}(a, b), \quad a, b > 0. \quad (7.5)$$

The joint posterior density of  $\beta, \sigma^2$  is

$$p(\beta, \sigma^2 | y) = p(\beta | y, \sigma^2) p(\sigma^2 | y)$$

so it suffices to compute  $p(\beta | y, \sigma^2)$  and  $p(\sigma^2 | y)$ .

In order to express the various conditional and marginal posteriors, we define a few variables to simplify notation:

$$\Omega := (M + X'X)^{-1} \quad (7.6)$$

$$\hat{\mu} = (M + X'X)^{-1}(X'X\hat{\beta} + M\beta_0) \quad (7.7)$$

$$Q := (\beta_0 - \hat{\beta})'(M^{-1} + (X'X)^{-1})^{-1}(\beta_0 - \hat{\beta}) \quad (7.8)$$

Then the conditional posterior  $p(\beta | \sigma^2, y)$  is normal with

$$\beta | \sigma^2, y \sim N(\hat{\mu}, \sigma^2 \Omega).$$

Note that if  $M = \lambda I$  and  $\beta_0 = 0$ , then  $\hat{\mu}$  is the ridge regression estimator  $\hat{\beta}_\lambda$ .

The marginal posterior for  $\sigma^2$  is:

$$\sigma^2 | y \sim \text{IG}\left(\frac{n}{2} + a, \frac{2b + s^2 + Q}{2}\right) \quad ? \quad (7.9)$$

Note that the posteriors for  $\beta$  and  $\sigma^2$  are of the same types as the prior (7.5), i.e. (7.5) is a conjugate prior.

Integrating out  $\sigma^2$  from  $p(\beta | \sigma^2, y)$  using the marginal posterior (7.9) yields the marginal posterior of  $\beta$ :

$$\beta | y \sim \mathcal{T}_{k+1}\left(n + 2a, \hat{\mu}, \frac{2b + s^2 + Q}{n + 2a} \Omega\right).$$

where the density of  $\mathcal{T}_p(\nu, \theta, \Sigma)$  is defined as

$$f(t | \nu, \theta, \Sigma) = \frac{\Gamma((\nu + p)/2)}{\Gamma(\nu/2) \sqrt{|\Sigma|} \nu \pi} \left[ 1 + \frac{(t - \theta)' \Sigma^{-1} (t - \theta)}{\nu} \right]^{-(\nu + p)/2}.$$

Hence  $\hat{\mu} = \mathbb{E}[\beta | y]$  is the posterior mean, which is the Bayes estimator associated to quadratic loss; it can be viewed as a weighted average between the maximum likelihood estimator and the prior mean.

**7.6. Black, Litterman, and Bayes.** The topic of portfolio optimization in the style of Black and Litterman (1991) and Black and Litterman (1992) seems to have generated more than its share of confusion over the years, as evidenced by articles with titles such as “*A demystification of the Black-Litterman model*” (Satchell and Scowcroft, 2000), etc. The method itself is often described as “Bayesian” but the original authors do not elaborate directly on connections with Bayesian statistics. The paper by Litterman and He (1999) contains many references to a “prior” but only one mention of a “posterior” without details, and no mention of a “likelihood.”

In the present lecture, we clarify the exact nature of the Bayesian statistical model to which Black-Litterman optimization corresponds, in terms of the prior, likelihood, and posterior. In the process we also lay out the full set of assumptions made, some of which are glossed over in other treatments.

Consider a view such as “the German equity market will outperform a capitalization-weighted basket of the rest of the European equity markets by 5%,” which is an example presented in Litterman and He (1999). Let  $\mathbf{p} \in \mathbb{R}^n$  denote a dollar-neutral portfolio which is long one unit of the DAX index, and short a one-unit basket which holds each of the other major European indices (UKX, CAC40, AEX, etc.) in proportion to their respective aggregate market capitalizations, so that  $\sum_i p_i = 0$ . Let  $q = 0.05$  in this example. This view may be equivalently expressed as

$$\mathbb{E}[\mathbf{p}' \mathbf{r}] = q \in \mathbb{R} \quad (7.10)$$

where  $\mathbf{r}$  is the random vector of asset returns over some subsequent interval, and  $q$  denotes the expected return, according to the view.

Suppose that the total number of such views is  $k$ , where each individual view is analogous to the above but with a different portfolio and corresponding expected return:

$$\mathbb{E}[\mathbf{p}_i' \mathbf{r}] = q_i, \quad i = 1 \dots k$$

then the portfolios  $\mathbf{p}_i$  are more conveniently arranged as rows of a matrix  $\mathbf{P}$ , and the statement of views becomes

$$\mathbb{E}[\mathbf{P}\mathbf{r}] = \mathbf{q} \text{ for } \mathbf{q} \in \mathbb{R}^k. \quad \mathbf{P} = \begin{matrix} & \rightarrow \mathbf{p}_1 \\ \rightarrow & \mathbf{p}_2 \\ \rightarrow & \mathbf{p}_3 \\ \rightarrow & \mathbf{p}_4 \\ \rightarrow & \mathbf{p}_5 \end{matrix} \quad (7.11)$$

The core idea of Black and Litterman (1991) is to treat the portfolio manager's views (7.11) as noisy observations which are useful for performing statistical inference concerning the parameters in some underlying model for  $\mathbf{r}$ . For example, if

$$\mathbf{r} \sim N(\boldsymbol{\theta}, \boldsymbol{\Sigma}) \quad (7.12)$$

with  $\boldsymbol{\Sigma}$  a known positive-definite  $n \times n$  matrix, then the views (7.11) can be recast as “observations” relevant for inference on the parameter  $\boldsymbol{\theta}$ . *try to get posterior of  $\theta$*

The practitioner must also specify a level of uncertainty or “error bar” for each view, which is assumed to be an independent source of noise from the volatility already accounted for in some underlying model for asset returns, such as (7.12). This is expressed as the following (more precise) restatement of (7.11):

$$\mathbb{E}(\mathbf{P}\mathbf{r}) = \mathbf{P}\boldsymbol{\theta} = \mathbf{q} + \boldsymbol{\epsilon}^{(v)}, \quad \boldsymbol{\epsilon}^{(v)} \sim N(0, \boldsymbol{\Omega}), \quad \boldsymbol{\Omega} = \text{diag}(\omega_1, \dots, \omega_k) \quad (7.13)$$

Eq. (7.13) specifies a Gaussian distribution and may be alternately written in the form of a likelihood:  *$q = \mathbf{P}\boldsymbol{\theta} + \Sigma^{(v)}$*

$$f(\mathbf{q} | \boldsymbol{\theta}) \propto \exp \left[ -\frac{1}{2}(\mathbf{P}\boldsymbol{\theta} - \mathbf{q})' \boldsymbol{\Omega}^{-1}(\mathbf{P}\boldsymbol{\theta} - \mathbf{q}) \right] \quad (7.14)$$

which is the standard normal likelihood for a multiple linear regression problem with dependent variable  $\mathbf{q}$  and design matrix  $\mathbf{P}$ .

Portfolio managers in this model specify *noisy, partial, indirect* information about  $\boldsymbol{\theta}$ , through their views. The information is *partial and indirect* because the views are on portfolio returns, i.e. linear transformations of returns, rather than on the asset returns directly. The information is *noisy*, with the noise modeled by  $\boldsymbol{\epsilon}^{(v)}$ , because the future is always uncertain.

A subjective, uncertain view about what will happen to a certain portfolio in the future is conceptually distinct from a noisy experimental observation such as an attempt to measure some physical constant with imperfect laboratory equipment. Nonetheless, for building intuition, we suggest thinking of a portfolio manager's forecast as an “observation of the future” in which the measuring device is a rather

murky and unreliable crystal ball. Only in this way is it analogous to the noisy measurements in experimental design which much of statistics is designed to model.

Quite generally, if any random variable  $r$  comes from a density  $p(r | \theta)$  with parameter  $\theta$ , and if one were given a set of noisy observations of realizations of  $r$ , then one could infer something about  $\theta$  by statistical inference. This would be the predicament of a physicist with a noisy measuring device, measuring a quantity that is itself random, and we suppose the physicist wants to know about the underlying data-generating process. Black and Litterman essentially say that the portfolio manager's view is, mathematically, no different from a noisy observation of a realization of (a linear transformation of) future returns.

As noted above, to perform statistical inference, observations alone are not sufficient; one needs to fully specify the statistical model, which includes a likelihood and a prior. In fact (7.13) specifies the likelihood as

$$f(\mathbf{q} | \boldsymbol{\theta}) \propto \exp \left[ -\frac{1}{2}(\mathbf{P}\boldsymbol{\theta} - \mathbf{q})'\boldsymbol{\Omega}^{-1}(\mathbf{P}\boldsymbol{\theta} - \mathbf{q}) \right]$$

which is the standard normal likelihood for a multiple linear regression problem with dependent variable  $\mathbf{q}$  and design matrix  $\mathbf{P}$ .

A feature of Bayesian statistics that is dissimilar from frequentist statistics is the ability to perform inference in data-scarce situations. In Bayesian statistics, even a single observation can lead to valid inferences for multi-parameter models due to the presence of a prior. In essence, when less information is available, more weight is given to the prior.

The classic regression problem has the number of variables much less than the number of observations, and is therefore identifiable. However, the need to perform inference in models with many more variables than observations also arises in many applications. Notably, this arises in the analysis of gene expression arrays, and is typically handled by Bayesian methods such as ridge and the lasso (Tibshirani, 1996).

In a Black-Litterman model with one single view, there is one observation and still  $n$  parameters to serve as the subjects for statistical inference:  $\boldsymbol{\theta} \in \mathbb{R}^n$  are the unobservable means of the asset returns. More generally, we may be presented with no views, one, or very many. When views are collected from many diverse portfolio managers or economists, they may contain internal contradictions; ie. it may be impossible that they all come true exactly. Bayesian regression is the ideal tool to deal with all such cases. Internal contradictions in the views simply mean that there is no exact (zero-residual) solution to the regression equations, which in fact is the typical situation in classic (identifiable) linear regression.

We have not yet specified the prior, but Black and Litterman were motivated by the guiding principle that, in the absence of any sort of information/views which



could constitute alpha over the benchmark, the optimization procedure should simply return the global CAPM equilibrium portfolio, with holdings denoted  $\mathbf{h}_{eq}$ . Hence in the absence of any views, and with prior mean equal to  $\mathbf{\Pi}$ , the investor's model of the world is that

$$\mathbf{r} \sim N(\theta, \Sigma), \quad \text{and} \quad \theta \sim N(\mathbf{\Pi}, \mathbf{C}) \quad (7.15)$$

for some covariance  $\mathbf{C}$  whose inverse represents the amount of precision in the prior.

For any portfolio  $\mathbf{p}$ , then, according to (7.15) we have

$$\mathbb{E}[\mathbf{p}'\mathbf{r}] = \mathbf{p}'\mathbf{\Pi} \quad \mathbb{V}[\mathbf{p}'\mathbf{r}] = \mathbf{p}'(\Sigma + \mathbf{C})\mathbf{p}. \quad (7.16)$$

$$\begin{aligned} \mathbb{E}[\mathbf{p}'\mathbf{r}] &= \mathbb{E}[\mathbb{E}[\mathbf{p}'\mathbf{r}|\theta, \Sigma]] = \mathbb{E}[\mathbf{p}'\theta] = \mathbf{p}'\mathbf{\Pi} \\ \text{Var}[\mathbf{p}'\mathbf{r}] &= \mathbb{E}[\text{Var}[\mathbf{p}'\mathbf{r}|\theta, \Sigma]] + \text{Var}[\mathbb{E}[\mathbf{p}'\mathbf{r}|\theta, \Sigma]] \\ &= \mathbb{E}[\mathbf{p}'\Sigma\mathbf{p}] + \text{Var}[\mathbf{p}'\theta] \\ &= \mathbf{p}'\Sigma\mathbf{p} + \mathbf{p}'\mathbf{C}\mathbf{p} \\ &= \mathbf{p}'(\Sigma + \mathbf{C})\mathbf{p} \end{aligned}$$

In fact we must make a choice whether to use the conditional or unconditional variance in optimization:

$$\mathbb{V}(\mathbf{r}|\theta) = \Sigma \quad \text{but} \quad \mathbb{V}(\mathbf{r}) = \Sigma + \mathbf{C}.$$

Since investors are presumably concerned with unconditional variance of wealth, the unconditional variance form is preferable, but when we take  $\mathbf{C} = \tau\Sigma$ , this distinction won't matter much (it amounts to a different risk-aversion constant).

Mean-variance optimization with (7.16) and with risk-aversion parameter  $\kappa > 0$ , leads to  $\kappa^{-1}(\Sigma + \mathbf{C})^{-1}\mathbf{\Pi}$  as the optimal portfolio; setting this equal to the CAPM equilibrium then gives

$$\theta \sim N(\mathbf{\Pi}, \mathbf{C}) \quad \mathbf{h}_{eq} = \kappa^{-1}(\Sigma + \mathbf{C})^{-1}\mathbf{\Pi}. \quad (7.17)$$

Any combination of  $\mathbf{\Pi}, \mathbf{C}$  satisfying (7.17) will lead to a model with the desired property – that the optimal portfolio with only the information given in the prior is the prescribed portfolio  $\mathbf{h}_{eq}$ . In particular, taking  $\mathbf{C} = \tau\Sigma$  with some arbitrary scalar  $\tau > 0$ , as did the original authors, leads to

$$\mathbf{\Pi} = \kappa(1 + \tau)\Sigma\mathbf{h}_{eq}$$

We thus have the normal likelihood (7.14) and the normal prior (7.15) which is a *conjugate prior* for that likelihood, meaning that the posterior is of the same family (ie. also normal in this example). A detailed discussion of conjugate priors is found in Robert (2007, Sec. 3.3). Let  $N(x; \mu, \sigma^2)$  denote the probability density of the normal, with  $\mu, \sigma$  as parameters and evaluated as a function of  $x$  for the density. Up to normalization factors,

$$\begin{aligned} \text{posterior} &\propto \underbrace{N(\mathbf{q}; \mathbf{P}\theta, \mathbf{\Omega})}_{\text{likelihood}} \times \underbrace{N(\theta; \mathbf{\Pi}, \mathbf{C})}_{\text{prior}} \\ &\propto \mathcal{P}(\mathbf{q}|\theta) \times \mathcal{P}(\theta) \end{aligned}$$

The normal density is of the form  $\exp[-\frac{1}{2}\text{quadratic}]$  so this will be easier to analyze if we take  $-2\log$  of both sides.

没有任何先验的  
view时，上述提到的P  
的最optimal的值

The negative log posterior is thus proportional to (neglecting terms that don't contain  $\theta$ ):

$$(P\theta - q)' \Omega^{-1} (P\theta - q) + (\theta - \Pi)' C^{-1} (\theta - \Pi) \quad (7.18)$$

$$= \theta' P' \Omega^{-1} P \theta - \theta' P' \Omega^{-1} q + q' \Omega^{-1} P \theta \quad (7.19)$$

$$+ \theta' C^{-1} \theta - \theta' C^{-1} \Pi - \Pi' C^{-1} \theta$$

$$= \theta' [P' \Omega^{-1} P + C^{-1}] \theta - 2(q' \Omega^{-1} P + \Pi' C^{-1}) \theta \quad (7.20)$$

一次项系数  $2\Pi' C^{-1} P + \Pi' C^{-1}$   
 $\mu = H^{-1} (q' \Omega^{-1} P + \Pi' C^{-1})$

Recall that, for  $H$  symmetric,  $\downarrow$  一次项系数对应的  $\frac{1}{2}$

$$\theta' H \theta - 2v' H \theta = (\theta - v)' H (\theta - v) - v' H v$$

Therefore if a multivariate normal random variable  $\theta$  has density  $p(\theta)$ ,

$$\begin{aligned} -2 \log p(\theta) &= \theta' H \theta - 2\eta' \theta + (\text{terms without } \theta) \\ \Rightarrow \mathbb{V}[\theta] &= H^{-1} \text{ and } \mathbb{E}\theta = H^{-1} \eta. \end{aligned}$$

For the quadratic term to match (7.20) we must have  $H = P' \Omega^{-1} P + C^{-1}$  and hence the posterior has mean

$$v = [P' \Omega^{-1} P + C^{-1}]^{-1} [P' \Omega^{-1} q + C^{-1} \Pi] \quad (7.21)$$

将括号展开会发现  $v$  是  $\theta$  和  $\Pi$  的线性组合。

and covariance

$$H^{-1} = [P' \Omega^{-1} P + C^{-1}]^{-1}. \quad (7.22)$$

Part of the beauty of this derivation is its simplicity: going from (7.18) to (7.22) requires just a few lines of algebra. Note that (7.22) is the posterior covariance of  $\theta$ , not the *a posteriori* covariance of  $r$ .

If asset returns are modeled using an elliptical distribution (and they are in this case), investors with any concave, increasing utility function will want to solve

$$h^* = \operatorname{argmax}_h \{ \mathbb{E}[h' r] - (\kappa/2) \mathbb{V}[h' r] \} \quad \mathbb{E}(h' r) = h' \mathbb{E}(r) = h' \mathbb{E}[\mathbb{E}(r|\theta)]$$

where  $\mathbb{E}[r]$  and  $\mathbb{V}[r]$  denote, respectively, the unconditional mean and covariance of  $r$  under the posterior. The unconditional covariance is a sum of variance due to parameter uncertainty, and variance due to the randomness in  $r$ . In other words,

$$\begin{aligned} \mathbb{V}[h' r] &= h' [P' \Omega^{-1} P + C^{-1}]^{-1} h + h' \Sigma h \\ &= h' \operatorname{var}(r) h = h' ( \mathbb{E}[\operatorname{var}(r|\theta, \Sigma)] + \operatorname{var}[\mathbb{E}(r|\theta, \Sigma)] ) h = h' (\Sigma + H^{-1}) h \end{aligned}$$

The optimal portfolio accounting for both types of variance is then

$$h^* = \kappa^{-1} [H^{-1} + \Sigma]^{-1} H^{-1} [P' \Omega^{-1} q + C^{-1} \Pi].$$

Just for fun, let's verify that this approaches  $h_{eq}$  when there aren't any views or, equivalently, when  $\Omega \rightarrow \infty$  or  $\Omega^{-1} \rightarrow 0$  and  $C = \tau \Sigma$ :

$$H^{-1} = [P' \Omega^{-1} P + C^{-1}]^{-1} = C = \tau \Sigma$$

这里的  $\Pi$  和  $C$  由  $\kappa'(\Sigma + C)^{-1} \Pi = h_{eq}$  得出。  
 此时 mean-variance optimization 的结果是  $h_{eq}$ 。  
 当  $q$  为 0 时  $\theta|q \sim N(Cv, H^{-1})$

mean variance optimization  
 $\left\{ \begin{aligned} & \text{find } h \text{ s.t. } h' v - \frac{1}{2} \kappa [h' (\Sigma + H^{-1}) h] \\ & \text{obj fn} = v - \kappa (\Sigma + H^{-1}) h = 0 \\ & \kappa (\Sigma + H^{-1}) h = v \\ & h^* = \frac{1}{\kappa} (\Sigma + H^{-1})^{-1} v \end{aligned} \right.$

$r \sim N(\theta, \Sigma), \theta \sim N(v, H^{-1})$

optimal prior portfolio.

optimal view portfolio

$\Rightarrow$  same combination of optimal view portfolio and optimal prior portfolio.

In this case

$$\begin{aligned} & \kappa^{-1}[\mathbf{H}^{-1} + \mathbf{\Sigma}]^{-1} \mathbf{H}^{-1} [\mathbf{P}' \mathbf{\Omega}^{-1} \mathbf{q} + \mathbf{C}^{-1} \mathbf{\Pi}] \\ &= [\kappa(1 + \tau) \mathbf{\Sigma}]^{-1} \tau \mathbf{\Sigma} [(\tau \mathbf{\Sigma})^{-1} \mathbf{\Pi}] \\ &= [\kappa(1 + \tau) \mathbf{\Sigma}]^{-1} \mathbf{\Pi} = \mathbf{h}_{eq} \end{aligned}$$

which verifies exactly the property we wanted to check.

## REFERENCES

- Bayes, Thomas (1763). “An essay towards solving a problem in the doctrine of chances.” In:
- Black, Fischer and Robert Litterman (1992). “Global portfolio optimization”. In: *Financial Analysts Journal*, pp. 28–43.
- Black, Fischer and Robert B Litterman (1991). “Asset allocation: combining investor views with market equilibrium”. In: *The Journal of Fixed Income* 1.2, pp. 7–18.
- Litterman, Robert and Guangliang He (1999). “The intuition behind Black-Litterman model portfolios”. In: *Goldman Sachs Investment Management Series*.
- Robert, Christian (2007). *The Bayesian choice: from decision-theoretic foundations to computational implementation*. Springer Science & Business Media.
- Satchell, Stephen and Alan Scowcroft (2000). “A demystification of the Black–Litterman model: Managing quantitative and traditional portfolio construction”. In: *Journal of Asset Management* 1.2, pp. 138–150.
- Tibshirani, Robert (1996). “Regression shrinkage and selection via the lasso”. In: *Journal of the Royal Statistical Society. Series B (Methodological)*, pp. 267–288.