Problem Formulation
$K = 1$ Recipe
Optimization Bias
Main Theorem
Appendix
Plots

# Optimization Bias in Covariance Estimation: Effect of Model Misspecification

Zhuoli Jin

Department of Statistics and Applied Probability

April 4, 2025

Problem Formulation
$K = 1$ Recipe
Optimization Bias
Main Theorem
Appendix
Plots

# Section 1

## Problem Formulation

Problem Formulation
$K = 1$ Recipe
Optimization Bias
Main Theorem
Appendix
Plots

## Overview

Given a $p \times n$ data matrix $Y$,

- The true (unknown) model is defined as

$$Y_{p \times n} = B_{p \times q} X_{q \times n} + \epsilon_{p \times n} \tag{1}$$

- The true covariance matrix $\Sigma$ is:

$$\Sigma = \mathscr{B} \Lambda \mathscr{B}^\top + \gamma^2 I, \tag{2}$$

- - $\Lambda$: $q \times q$ diagonal matrix, eigenvalues of $BB^\top$
  - $\mathscr{B}$: $p \times q$ matrix of corresponding eigenvectors
  - $\gamma$: $E(\epsilon^\top \epsilon) = \gamma^2 I$

We don't know $q \Rightarrow$ We choose a $K$.

Problem Formulation
$K = 1$ Recipe
Optimization Bias
Main Theorem
Appendix
Plots

## $K = 1$

Assume the true number of factors $q > 1$, and we choose the estimate $K = 1$. The population covariance matrix that we _believe_ is

$$\Sigma^* = \sigma_p^2 bb^\top + \gamma^2 I, \qquad (3)$$

where $\sigma_p^2$ is the largest eigenvalue of $bb^\top$. This motivates the estimated covariance matrix as

$$\hat{\Sigma} = s^2 hh^\top + \hat{\gamma}^2 I. \qquad (4)$$

Note: we use $^*$ to denote our _belief_ and $\hat{\ }$ to denote the _estimate_.

Problem Formulation
$K = 1$ Recipe
Optimization Bias
Main Theorem
Appendix
Plots

## Notations & Covariance Matrix Models

| | |
|---|---|
| $q$ | True number of spikes |
| $K$ | Estimated number of spikes |
| $\mathscr{B}$ | Eigenvectors of $BB^\top$ |
| $b$ | First column of $\mathscr{B}$, normalized to length 1 |
| $h$ | Leading eigenvector of $\hat{\Sigma} = YY^\top/n$ |
| $s_{i,p}^2$ | The $i$-th eigenvalue of $\hat{\Sigma} = YY^\top/n$ |
| $\lambda_i^2$ | The $i-$th eigenvalue of $L = Y^\top Y/p$ |

### Table: Notations

| | |
|---|---|
| $\Sigma$ | True (unknown) population covariance matrix |
| $\Sigma^*$ | Covariance matrix when we choose $K = 1$; $\Sigma^* = \sigma_p^2 bb^\top + \gamma^2 I$ |
| $\hat{\Sigma}$ | Sample covariance matrix based on $K = 1$; $\hat{\Sigma} = s^2 hh^\top + \hat{\gamma}^2 I$ |

### Table: Comparison Table

# Section 2

## $K = 1$ Recipe

Problem Formulation
$K = 1$ Recipe
Optimization Bias
Main Theorem
Appendix
Plots

## $h_\sharp$ Construction

Recipe:

$$h_\sharp = \frac{1}{D}(\psi^2 h + N z_{\perp h}), \tag{5}$$

where

$$z_{\perp h} = \frac{z - z_h}{|z - z_h|}$$

with

$$z_h = \langle h, z \rangle h$$

.

- $z$: the direction that we want to study
- $z_h$: the projection of $z$ onto $h$
- $z_{\perp h}$: the unit vector in the direction of the component of $z$ that's orthogonal to $h$
- $D$: a normalizing constant
- $\psi^2$: signal-to-noise ratio

Problem Formulation
K = 1 Recipe
Optimization Bias
Main Theorem
Appendix
Plots

## Details in $h_\sharp$

Define

$$\psi^2 = \frac{s_{1,p}^2 - l_p^2}{s_{1,p}^2}, \text{ with } l_p^2 \text{ be average non-zero bulk eigenvalues.}$$

$$N = \frac{\langle h, z \rangle - \psi^2 \langle h, z \rangle}{\sqrt{1 - \langle h, z \rangle^2}}, \tag{6}$$

$$D = \sqrt{\psi^4 + N^2}$$

- $\psi^2$: signal-to-noise ratio. $\psi^2$ is high – put more weight on $h$
- $N$: the correction strength in the direction orthogonal to $h$, scaled by the noise level $(1 - \psi^2)$ and alignment between $h$ and $z$. $N$ is high – need to pull back by adding more weight on $z_{\perp h}$.

# Section 3

## Optimization Bias

Problem Formulation
$K = 1$ Recipe
**Optimization Bias**
Main Theorem
Appendix
Plots

## $\mathscr{E}_p(h)$

For $z \in \mathbb{R}^p$ with $|z| = 1$, the <u>*true*</u> quadratic optimization function is defined as:

$$\mathscr{E}_p(h) = \frac{\mathscr{B}^\top z - \mathscr{B}^\top \mathscr{H} \mathscr{H}^\top z}{\sqrt{1 - \langle \mathscr{H}, z \rangle^2}} \qquad (7)$$

Connecting to minimum variance problem:

The expected out-of-sample variance $V_p^2 = \langle \hat{w}, \Sigma \hat{w} \rangle$ can be written as

$$V_p^2 = \frac{|\Lambda_p \mathscr{E}_p(\mathscr{H})|^2}{p|z - z_{\mathscr{H}}|^2} + O(1/p), \qquad (8)$$

so $V_p^2 \to 0 \iff |\mathscr{E}_p(\mathscr{H})| \to 0$.

Problem Formulation
$K = 1$ Recipe
Optimization Bias
Main Theorem
Appendix
Plots

# $\mathscr{E}_p^*(h)$

We _believe_ the optimization bias is

$$\mathscr{E}_p^*(h) = \frac{\langle b, z \rangle - \langle b, h \rangle \langle h, z \rangle}{\sqrt{1 - \langle h, z \rangle^2}}. \tag{9}$$
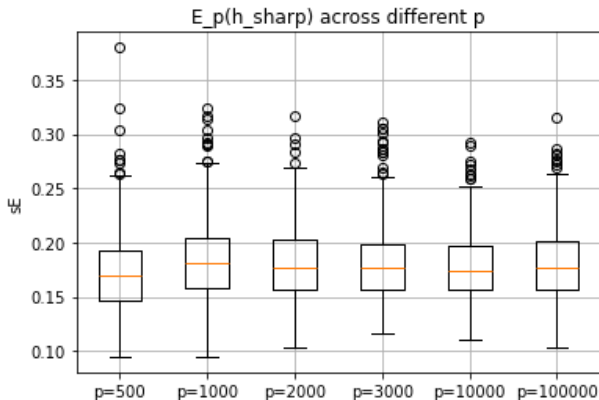
Note it is the first component of the true $\mathscr{E}_p(h)$ (by writing $B = (\beta_1, \beta_2, \cdots, \beta_q)$ and $b = \beta_1/|\beta_1|$).

When $q = K = 1$, $\mathscr{E}_p^*(h) \to 0$.

Problem Formulation
$K = 1$ Recipe
Optimization Bias
Main Theorem
Appendix
Plots

# $|E_p^*(h)|$ v.s. $p$ $(q = 7)$

Problem Formulation
$K = 1$ Recipe
Optimization Bias
Main Theorem
Appendix
Plots

# $|E_p^*(h_\sharp)|$ v.s. $p$ $(q = 7)$

# Section 4

## Main Theorem

Problem Formulation
$K = 1$ Recipe
Optimization Bias
**Main Theorem**
Appendix
Plots

## Theorem

### Theorem

*For $\psi^2 = \frac{s_p^2 - l_p^2}{s_p^2}$, where $l_p^2$ is the average non-zero eigenvalues, we have the following results:*

1. $\mathscr{E}_p^*(h_\sharp) = \frac{\psi^2 \langle b, z \rangle - \langle h, b \rangle \langle h, z \rangle}{\sqrt{\psi^4 + (1 - 2\psi^2) \langle h, z \rangle^2}}$.

2. *When $K = 1 = q$, $\mathscr{E}_p^*(h_\sharp) \to 0$*

3. *When $K = 1 < q$, $\mathscr{E}_p^*(h_\sharp) \not\to 0$.*

Problem Formulation
$K = 1$ Recipe
Optimization Bias
Main Theorem
Appendix
Plots

# $|E_p^*|$ v.s. $p$ ($q = 7$)

Problem Formulation
K = 1 Recipe
Optimization Bias
Main Theorem
Appendix
Plots

# $|E_p^*|$ v.s. $q$ I



Figure: 1 spiked factor volatility in $B$ construction. $B$: [market style block] $\times$ factor return matrix

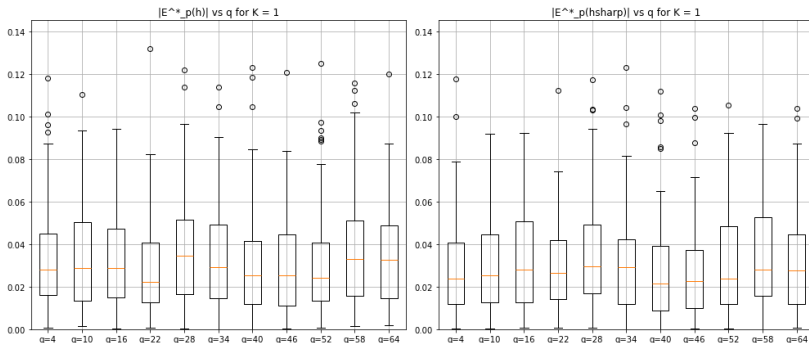Problem Formulation
$K = 1$ Recipe
Optimization Bias
**Main Theorem**
Appendix
Plots

# $|E_p^*|$ v.s. $q$ II



Figure: 2 spiked factor volatilities in $B$ construction. $B$: [market style block] $\times$ factor return matrix

Problem Formulation
$K = 1$ Recipe
Optimization Bias
**Main Theorem**
Appendix
Plots

# $|E_p^*|$ v.s. $q$ for Different Factor Volatilities I

Note: the first two diagonal elements of factor vol matrix is 16 and 8, the remains are randomly drawed integers between 1 and $M$.
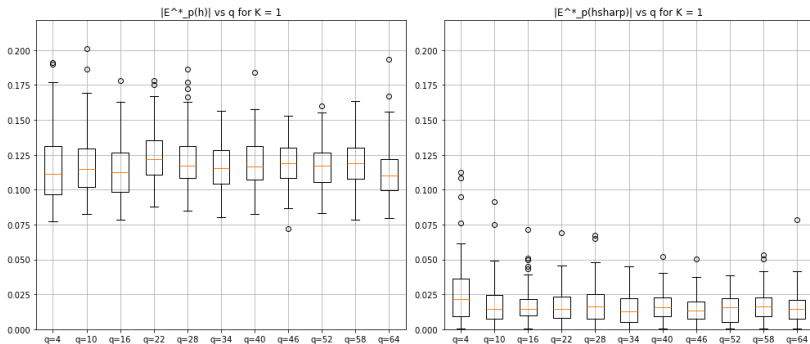


Figure: $M = 25$

Problem Formulation
$K = 1$ Recipe
Optimization Bias
**Main Theorem**
Appendix
Plots
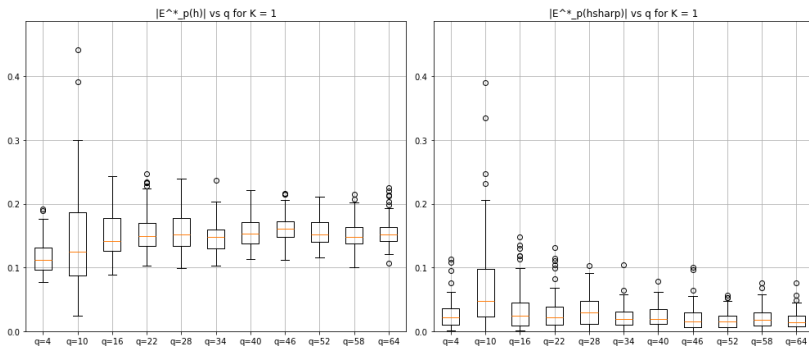
# $|E_p^*|$ v.s. $q$ for Different Factor Volatilities II



Figure: $M = 50$
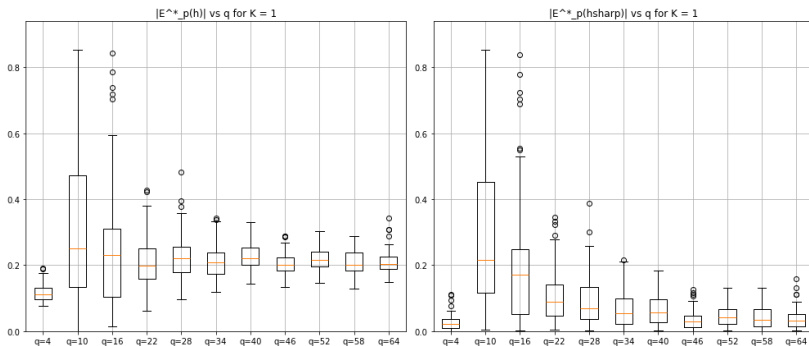
# $|E_p^*|$ v.s. $q$ for Different Factor Volatilities III



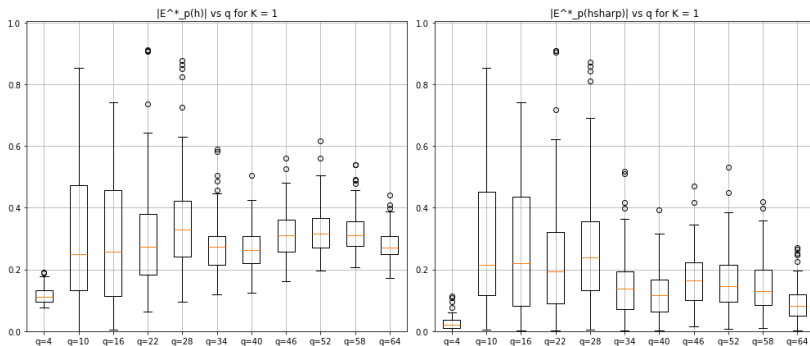Figure: $M = 75$

Problem Formulation
K = 1 Recipe
Optimization Bias
Main Theorem
Appendix
Plots

# $|E_p^*|$ v.s. $q$ for Different Factor Volatilities IV



Figure: $M = 100$

# Section 5

## Appendix

Problem Formulation
$K = 1$ Recipe
Optimization Bias
Main Theorem
**Appendix**
Plots

## $q = 1$ Asymptotics

If $K = q = 1$, we have

1. $\lim_{p\uparrow\infty} |\langle h, b \rangle^2 - \psi^2| = 0$ for $\psi^2 = 1 - \kappa_p^2 s_{1,p}^{-2}$, where
   $\kappa_p^2 = \frac{\sum_{j>q} s_{j,p}^2}{n-q}$

2. $\lim_{p\uparrow\infty} |\langle h, z \rangle - \langle h, b \rangle\langle b, z \rangle| = 0$

As $p \to \infty$,

$$\begin{aligned}
\langle b, z \rangle - \langle b, h \rangle\langle h, z \rangle &= \langle b, z \rangle - \langle b, h \rangle\langle h, b \rangle\langle b, z \rangle \\
&= \langle b, z \rangle(1 - \langle b, h \rangle\langle h, b \rangle) \neq 0
\end{aligned} \quad (10)$$

Problem Formulation
$K = 1$ Recipe
Optimization Bias
Main Theorem
**Appendix**
Plots

# $\mathscr{E}_p^*(h_\sharp)$

Recall $\mathscr{E}_p^*(h_\sharp) = \mathscr{E}_p^*(h_z t_\sharp)$. Define $\tilde{\mathscr{E}}_p(\cdot) : t \mapsto \mathscr{E}_p^*(h_z t)$. Claim: $t \mapsto \tilde{\mathscr{E}}_p(t)$ is continuous in $\mathbb{R}$.

1. $\langle h_z t, z \rangle < 1$
2. $\langle b, z \rangle - \langle b, h_z t \rangle \langle h_z t, z \rangle = \langle b, z \rangle - t^2 \langle b, h_z \rangle \langle h_z, z \rangle$ is continuous w.r.t $t$.

Now we have for $1 = K = q$:

$$\left. \begin{array}{l} \tilde{\mathscr{E}}_p(t) \text{ is continuous} \\ \tilde{\mathscr{E}}_p(t_\sharp) = \tilde{\mathscr{E}}_p(t_\sharp - t_* + t_*) \\ \tilde{\mathscr{E}}_p(t_*) = \mathscr{E}_p^*(h_z t_*) = 0 \\ |t_\sharp - t_*| \to 0 \end{array} \right\} \Rightarrow \mathscr{E}_p^*(h_\sharp) \to 0. \qquad (11)$$

When $1 = K < q$, $|t_\sharp - t_*| \nrightarrow 0$.

Problem Formulation
$K = 1$ Recipe
Optimization Bias
Main Theorem
**Appendix**
Plots

# $\mathscr{E}_p^*(h_\sharp)$ (Continued)

Let $D = \sqrt{\psi^4 + \frac{(\langle h, z\rangle - \psi^2\langle h, z\rangle)^2}{1 - \langle h, z\rangle^2}}$ and $N = \frac{\langle h, z\rangle - \psi^2\langle h, z\rangle}{\sqrt{1 - \langle h, z\rangle^2}}$. Note that

1. $\langle h_\sharp, z\rangle = \frac{\langle h, z\rangle}{D}$

2. $\langle h_\sharp, b\rangle = \frac{\psi^2\langle h, b\rangle + (1 - \psi^2)\langle h, z\rangle\langle z, b\rangle - \langle h, b\rangle\langle h, z\rangle^2}{D(1 - \langle h, z\rangle^2)}$

we have

$$
\begin{aligned}
\mathscr{E}_p^*(h_\sharp) &= \frac{D\langle b, z\rangle}{\sqrt{D^2 - \langle h, z\rangle^2}} - \frac{\psi^2\langle h, b\rangle\langle h, z\rangle + (1 - \psi^2)\langle h, z\rangle^2\langle b, z\rangle - \langle h, z\rangle^3\langle h, b\rangle}{D(1 - \langle h, z\rangle^2)\sqrt{D^2 - \langle h, z\rangle^2}} \\
&= \frac{(\psi^2 - \langle h, z\rangle^2)(\psi^2\langle b, z\rangle - \langle h, b\rangle\langle h, z\rangle)}{D(1 - \langle h, z\rangle^2)\sqrt{D^2 - \langle h, z\rangle^2}} \\
&= \frac{\psi^2\langle b, z\rangle - \langle h, b\rangle\langle h, z\rangle}{\sqrt{\psi^4 + (1 - 2\psi^2)\langle h, z\rangle^2}}
\end{aligned}
$$

(12)

Problem Formulation
$K = 1$ Recipe
Optimization Bias
Main Theorem
**Appendix**
Plots

# $\phi^2$ I

Denote the columns of $B$ be $\beta_1, \beta_2, \cdots, \beta_q$, with each $\beta_i \in \mathbb{R}^p$. Define

$$
\begin{aligned}
D_1 &= \begin{pmatrix} (\beta_1^\top b)(b^\top \beta_1) - \beta_1^\top \beta_1 & \cdots & (\beta_1^\top b)(b^\top \beta_q) - \beta_1^\top \beta_q \\ \vdots & \ddots & \vdots \\ (\beta_q^\top b)(b^\top \beta_1) - \beta_q^\top \beta_1 & \cdots & (\beta_q^\top b)(b^\top \beta_q) - \beta_q^\top \beta_q \end{pmatrix}, \\
D_2 &= \begin{pmatrix} (\beta_1^\top b)\beta_1 - \beta_1^\top & \cdots & (\beta_1^\top b)\beta_q - \beta_1^\top \\ \vdots & \ddots & \vdots \\ (\beta_q^\top b)\beta_1 - \beta_q^\top & \cdots & (\beta_q^\top b)\beta_q - \beta_q^\top \end{pmatrix}, \\
D_3 &= \begin{pmatrix} (bb^\top - 1)\beta_1 & \cdots & (bb^\top - 1)\beta_q \end{pmatrix}
\end{aligned}
\tag{13}
$$

And

$$
M_p = \frac{\epsilon^\top bb^\top \epsilon}{p}, \ \Gamma_p = \frac{\epsilon^\top \epsilon}{p}, \ N_p = \frac{X^\top D_1 X + \epsilon^\top D_2 X + X^\top D_3 \epsilon}{p}.
\tag{14}
$$

Problem Formulation
$K = 1$ Recipe
Optimization Bias
Main Theorem
**Appendix**
Plots

# $\phi^2$ II

Let $\omega = Y^\top h/(s_{1,p}\sqrt{n})$ denote the leading eigenvector of the dual covariance matrix $L$ with eigenvalue $\lambda_1^2$, and $\mathcal{W} = \omega\sqrt{p/(ns_{1,p}^2)}$, then we have

---

### Theorem

1. $\langle h, b \rangle^2 = \omega^\top(L + M_p - \Gamma_p + N_p)\omega/\lambda_1^2 = 1 + \omega^\top(M_p - \Gamma_p + N_p)\omega/\lambda_1^2$

2. Let $\phi^2 = \psi^2 + \mathcal{W}^\top(N_p + \frac{n\kappa_p^2}{p}I - \Gamma_p)\mathcal{W}$, then $|\langle h, b \rangle^2 - \phi^2| \to 0$.

---

- When true $q = 1$, the signal space effectively reduces to the single axis spanned by $b$. The largest eigenvalue is well separated from any noise directions, so the sample covariance matrix naturally aligns $h$ with $b$.

- When $q > 1$, the signal space becomes multi-dimensional. The estimated $h$ will always pick up contributions from other spikes. Consequently, $\langle h, b \rangle^2$ cannot match the single-spike formula $\psi^2$.

Problem Formulation
$K = 1$ Recipe
Optimization Bias
Main Theorem
**Appendix**
Plots

# $1 = K < q$ Asymptotics

Denote $Z_p = bb^\top BX + bb^\top \epsilon$, we have

$$
\begin{aligned}
|\langle h, z \rangle - \langle h, b \rangle \langle b, z \rangle| &= |h^\top z - (bb^\top h)^\top z| \\
&= \frac{1}{\sqrt{p}} \left| \mathcal{W}^\top X^\top B^\top z + \mathcal{W}^\top \epsilon^\top z - \mathcal{W}^\top Z_p^\top z \right| \\
&= \frac{1}{\sqrt{p}} \left| \mathcal{W}^\top \epsilon^\top (I - bb^\top) z + \mathcal{W}^\top X^\top B^\top (I - bb^\top) z \right|
\end{aligned}
\tag{15}
$$

Since $\overline{\lim}_{p \to \infty} |\mathcal{W}^\top| < \infty$ and $|\epsilon^\top (I - bb^\top) z| / \sqrt{p} \to 0$, we have

---

**Lemma**

1. $|\langle h, z \rangle - \langle h, b \rangle \langle b, z \rangle| \sim \frac{|\mathcal{W}^\top X^\top B^\top (I - bb^\top) z|}{\sqrt{p}}$

2. $\left| \langle h, b \rangle \mathscr{E}_p^*(h) - \frac{\langle h, z \rangle - \phi^2 \langle h, z \rangle}{\sqrt{1 - \langle h, z \rangle^2}} \right| \sim \frac{|\mathcal{W}^\top X^\top B^\top (I - bb^\top) z|}{\sqrt{p(1 - \langle h, z \rangle^2)}} > 0.$
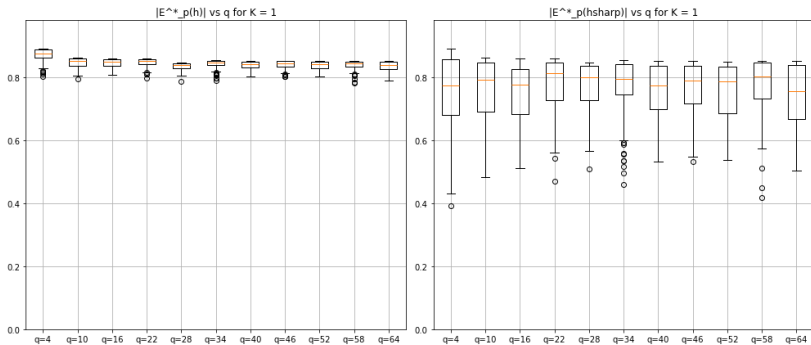
---

# Section 6

## Plots

Problem Formulation
$K = 1$ Recipe
Optimization Bias
Main Theorem
Appendix
Plots

# $\left|E_p^*\right|$ v.s. $q$ I



Figure: Factor Vol $= 1$

Problem Formulation
$K = 1$ Recipe
Optimization Bias
Main Theorem
Appendix
Plots

# $\left|E_p^*\right|$ v.s. $q$ II



Figure: Factor Vol = 20

Problem Formulation
$K = 1$ Recipe
Optimization Bias
Main Theorem
Appendix
**Plots**

# $\psi^2$ v.s. $p$

Problem Formulation
$K = 1$ Recipe
Optimization Bias
Main Theorem
Appendix
**Plots**

# $|\langle h, b \rangle|$ v.s. $p$

Problem Formulation
$K = 1$ Recipe
Optimization Bias
Main Theorem
Appendix
Plots

# $\langle h, z \rangle^2$ v.s. $p$



|<h,z>| across different p